

**FROM SPATIO-TEMPORAL DATA TO A WEIGHTED AND
LAGGED NETWORK BETWEEN FUNCTIONAL DOMAINS:
APPLICATIONS IN CLIMATE AND NEUROSCIENCE**

A Thesis
Presented to
The Academic Faculty

by

Ilias Fountalis

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
School of Computer Science

Georgia Institute of Technology
May 2016

Copyright © 2016 by Ilias Fountalis

**FROM SPATIO-TEMPORAL DATA TO A WEIGHTED AND
LAGGED NETWORK BETWEEN FUNCTIONAL DOMAINS:
APPLICATIONS IN CLIMATE AND NEUROSCIENCE**

Approved by:

Professor Constantine Dovrolis, Advisor
School of Computer Science
Georgia Tech

Professor Mostafa H. Ammar
School of Computer Science
Georgia Tech

Assistant Professor Bistra Dilkina
School of Computational Science and
Engineering
Georgia Tech

Professor Annalisa Bracco
School of Earth and Atmospheric
Sciences
Georgia Tech

Professor Athanasios Nenes
School of Earth and Atmospheric
Sciences
Georgia Tech

Associate Professor Shella Keilholz
Wallace H. Coulter Department of
Biomedical Engineering
Georgia Tech

Date Approved: 30 March 2016

To my family.

To my friends.

To Saamer.

ACKNOWLEDGEMENTS

I joined Georgia Tech at the age of twenty five and after five and a half years I am at a point where an interesting new chapter in my life begins. I was lucky enough to have Constantine's support and guidance throughout this time. His knack on finding interesting research problems made this thesis possible. His attention to detail and his continuous encouragement for me to understand in depth every single aspect of my research, changed me as a scientist and as a person. I am proud to consider Constantine as my mentor and as a valuable friend.

I would also like to thank Annalisa Bracco, which I consider as a co-advisor on this interesting research journey. Annalisa helped me understand concepts that (being a computer scientist at heart) I was totally unfamiliar with. Special thanks go to Athanasios Nenes with whom we exchanged interesting research ideas and whose advice in difficult times helped me stay on the right track. I would also like to thank Bistra Dilkina and Shella Keilholz. The last part of my thesis would not have been possible without their contribution.

I would also like to thank the NTG faculty for their continuous support and all the people at the NTG lab for their company and friendship. Special thanks go to Demetris Antoniadis and Anirudh Ramachandran. I would like to thank my mother Anna and my father Dimos for their unconditional love and for supporting me in all my choices. Last but not least, I would like to thank my Atlanta friends. I would never have managed to go through this journey without their help and support.

Contents

ACKNOWLEDGEMENTS	iv
LIST OF TABLES	viii
LIST OF FIGURES	ix
SUMMARY	xvi
I INTRODUCTION	1
1.1 Dimensionality reduction methods for spatio-temporal data	2
1.2 A framework for the analysis of spatio-temporal systems	3
1.2.1 <i>geo-Cluster</i>	3
1.2.2 δ -MAPS	4
1.3 Organization	5
II RELATED WORK	6
2.1 Synthetic Data Generation	6
2.2 Network Based Methods	7
2.3 Dimensionality Reduction methods	8
2.3.1 Principal Component Analysis	10
2.3.2 Independent Component Analysis	11
2.3.3 Clustering based methods	11
2.3.4 Community Detection Methods	14
III GEO-CLUSTER: SPATIO-TEMPORAL NETWORK ANALYSIS FOR STUDY- ING CLIMATE PATTERNS	17
3.1 Introduction	17
3.2 Data sets	20
3.3 Climate network construction	20
3.3.1 Cell-level network	21
3.3.2 Identification of climate areas	22
3.3.3 Links between areas	24

3.4	Network metrics	27
3.5	Robustness analysis	34
3.5.1	Robustness to additive white Gaussian noise	35
3.5.2	Robustness to the resolution of the input data set	36
3.5.3	Robustness to the selection of τ	36
3.5.4	Robustness to the selection of the correlation metric	37
3.6	Applications	41
3.6.1	Comparison of SST networks	41
3.6.2	Network changes over time	47
3.6.3	Comparison of precipitation networks	48
3.6.4	<i>Regression between networks</i>	51
3.6.5	CMIP5 SST networks	53
3.7	Discussion and Conclusions	58
3.8	Selection of threshold τ	60
3.9	Pseudocode of area identification algorithm	61
IV	ENSO IN CMIP5 SIMULATIONS: NETWORK CONNECTIVITY FROM THE RECENT PAST TO THE TWENTY-THIRD CENTURY	64
4.1	Introduction	64
4.2	Climate Network Inference	66
4.3	Results	71
4.3.1	CMIP5 Models and Observational Datasets	71
4.3.2	The Historical Experiments: 1956-2005	72
4.3.3	The RCP8.5 Experiments: 2051-2100	76
4.3.4	The ECP8.5 Experiments: 2101 - 2300	79
4.4	Discussion	88
4.5	Supplementary strength and link maps	91
4.6	Advantages of using a complete weighted cell-level network	104
V	δ-MAPS: FROM SPATIO-TEMPORAL DATA TO A WEIGHTED AND LAGGED NETWORK BETWEEN FUNCTIONAL DOMAINS	107

5.1	Introduction	107
5.2	Related Work	109
5.3	δ -MAPS	110
5.3.1	Functional domains	111
5.3.2	The domain network	114
5.4	Illustration - Comparisons	118
5.5	Application in Climate Science	119
5.6	Application in fMRI data	123
5.7	Discussion	126
5.8	Identifying the largest domain is NP-complete	127
5.9	Heuristic for the selection of δ	128
5.10	δ -MAPS pseudocode	129
VI	CONCLUSIONS & FUTURE WORK	132
6.1	Conclusions	132
6.2	Future Work	134
	REFERENCES	137

List of Tables

1	Synthetic domain generation parameters.	7
2	D_{sd} and ARI from HadISST (1979-2005) to reanalyses, GISS-E2H and HadCM3, and corresponding noise-to-signal ratios γ	58
3	List of models analyzed and global mean trends in sea surface temperature and rainfall over 1956-2005 and 2051-2100. The number of ensemble members considered during the historical period (1956-2005) is indicated for each model. In parenthesis the number of members with projections to 2100 under the RCP8.5 scenario. X indicates that the model has one member continuing to 2300. Boreal winter (December to February) global mean trends are averaged over all ensemble members (\pm denotes the maximum deviation between ensemble members)	73
4	Projected global mean trends in sea surface temperature and rainfall from 2101 to 2300. Trends are calculated over 50-year long consecutive intervals for the models with one member extending to 2300 and for boreal winter (December to February). Precipitation trends are in parenthesis	83

List of Figures

1	A: The five ground-truth domains. Adjacent domains have different colors, overlapping regions shown in black, and the core of each domain is in blue. The three constructed edges are shown in gray lines. B: The homogeneity field $\hat{r}_K(i)$ at each cell. The identified seeds are shown in blue. C: The inferred domains: adjacent domains have different colors and overlaps are shown in black. D: The inferred domain-level network: the color map refers to the edge correlation. The lag associated with each edge is also shown. E,F,G: The first three EOF (PCA) components. The variance explained by each component is shown at the top of each figure. H,I: The two ICA components. J,K: K-means clustering. L: The second hierarchical level of community structure as identified by OSLOM: each community has a distinct color and overlaps are shown in black.	9
2	Empirical Cumulative Distribution Functions (CDF) of correlations for the HadISST reanalysis during the 1950-1976 and 1979-2005 periods, and for ERSST-V3 and NCEP data during the 1979-2005 period	22
3	An example of the area identification algorithm. (a) 12-cell synthetic grid. (b) The correlation matrix between cells (given as input). (c) The area expansion process for a given $\tau=0.4$. Cells shown in red are selected to join the area (denoted by A_k). Cells 1, 4, 9 and 12 will not join A_k since they do not satisfy the τ constraint in Eq.2	24
4	Identified areas in the HadISST 1979-2005 data set ($\tau=0.496$). (a) The 176 areas identified by Part-1 of the area identification algorithm. (b) The 74 “merged” areas after the execution of Part-2. (c) The CDF of area sizes (in number of cells) before and after the merging process	25
5	The relation between area size and standard deviation of the area’s cumulative anomaly ($R^2 = 0.88$) for the HadISST reanalysis during the 1979-2005 period; $\tau=0.496$	26
6	CDF of the absolute correlation between area cumulative anomalies for the HadISST reanalysis during the 1950-1976 and 1979-2005 periods, and for ERSST-V3 and NCEP during the 1979-2005 period	27
7	Link maps for two areas related to (a) ENSO and (b) the equatorial Indian Ocean in the HadISST 1979-2005 network ($\tau=0.496$). The color scale represents the weight of the link between the area shown in black and every other area in this SST network	28
8	Strength maps for two different time periods using the HadISST data set. (a) 1950-1976 network, strength of ENSO area: 20.1×10^4 ; (b) 1979-2005 network, strength of ENSO area: 18.8×10^4	29

9	Color maps depicting the <i>top-5 order cores</i> for the (a) HadISST 1950-1976, and (b) HadISST 1979-2005 networks	30
10	(a) Distribution of ranked area strengths for two networks constructed using the HadISST data set over the periods 1950-1976 and 1979-2005, respectively. (b) Distance $D_{sd}(N, N_\gamma)$ and $ARI(N, N_\gamma)$ between the HadISST 1979-2005 network and networks constructed after the addition of white Gaussian noise in the same data set	33
11	Strength maps for two perturbations of the HadISST 1979-2005 data set using white Gaussian noise. (a) $\gamma=0.05$, strength of ENSO area: 18.0×10^4 . (b) $\gamma=0.10$, strength of ENSO area: 19.1×10^4	35
12	Strength maps for the HadISST 1979-2005 network at three different resolutions. (a) Low resolution network, ($4^\circ lat \times 4^\circ lon$), strength of ENSO area: 18.2×10^4 . (b) Default resolution network, ($2^\circ lat \times 2.5^\circ lon$), strength of ENSO area: 18.8×10^4 . (c) High resolution network, ($1^\circ lat \times 2^\circ lon$), strength of ENSO area: 18.2×10^4	38
13	(a) Distance D_{sd} and (b) ARI from the original HadISST 1979-2005 network (marked with an asterisk in the x-axis, $\tau=0.496$) to networks constructed with different values of τ . The black horizontal lines correspond to the distance $D_{sd}(N, N_\gamma)$ and $ARI(N, N_\gamma)$	39
14	Strength maps for the HadISST 1979-2005 network using two values of the parameter τ . The “default” value is $\tau=0.496$, corresponding to $\alpha=.1\%$ (see Section 3.8). (a) $\tau=0.45$, strength of ENSO area: 18.7×10^4 . (b) $\tau=0.55$, strength of ENSO area: 18.6×10^4	40
15	Strength map for the HadISST 1979-2005 network using Spearman’s correlation; strength of ENSO area: 18.5×10^4	40
16	Pearson correlation maps between the SST anomaly time series in all pairs of three reanalyses data sets over the 1979-2005 period in boreal winter (DJF). Correlations between (a) HadISST and ERSST-V3; (b) HadISST and NCEP; (c) NCEP and ERSST-V3	43
17	Strength maps for networks constructed based on (a) HadISST (ENSO area strength 18.8×10^4); (b) ERSST-V3 (ENSO area strength 17.6×10^4); (c) NCEP (ENSO area strength 21.0×10^4). In all networks the period considered is 1979-2005	44
18	<i>Top-5 order cores</i> in (a) HadISST; (b) ERSST-V3; (c) NCEP. The period considered is 1979-2005 in all cases	45
19	Links between the ENSO-like area shown in black and all other areas in the three reanalyses. (a) HadISST, (b) ERSST-V3 and (c) NCEP networks .	46

20	Links for the HadISST network over 1950 - 1976 from the (a) ENSO-related area, and (b) the equatorial Indian Ocean area (in black in the two panels)	49
21	Precipitation networks. Area strength map in (a) CMAP (equatorial Pacific area strength 49.4×10^4), and (b) ERA-Interim (equatorial area strength 41.0×10^4)	50
22	<i>Top-5 order cores</i> in (a) CMAP, and (b) ERA-Interim	50
23	Link maps from the strongest area (in black) for the two precipitation re-analysis data sets. (a) CMAP; (b) ERA Interim	52
24	Link maps from the ENSO-like area in HadISST data set to all areas in the CMAP data set, considering the 1979-2005 period. Values greater than $ 1 \times 10^4 $ are saturated	53
25	Strength maps for two members of the GISS-E2H and HadCM3 “historical” ensemble. (a) GISS-E2H run 1 (ENSO area strength 9.8×10^4); (b) GISS-E2H run 2 (ENSO area strength 10.0×10^4); (c) HadCM3 run 1 (ENSO area strength 23.3×10^4) and (d) HadCM3 run 2 (ENSO area strength 16.9×10^4)	55
26	<i>Top-5 order cores</i> identified in the SST anomaly networks for (a-b) two GISS-E2H ensemble members and (c-d) two HadCM3 integrations	56
27	Link maps from the ENSO-like area in the (a-b) GISS-E2H and (c-d) HadCM3 models	57
30	Trend anomaly maps for boreal winter in the recent past and near future. Anomalies are computed by removing the global mean trend calculated over the months of December to February and indicated in Table 3 from each grid cell. + and • indicate agreement in more than 90% and 70% of models in the sign of the trend anomaly slope. (a) HadISST. (b) ERA40+Interim. (c) Sea surface temperature (SST) averaged across models in the historical period (1956-2005). (d) As in (c) but for rainfall. The units are C° /year for SST and (mm/day)/year for precipitation	74
31	Metric D versus ARI for climate networks during the historical period 1956-2005. (a) Sea surface temperature; reference network HadISST. (b) Precipitation; reference network ERA40+Interim. Three levels of noise-to-signal ratios γ are also indicated	77

32	Strength maps of sea surface temperature for HadISST and three sample models (top rows), and of precipitation for ERA40+Interim and the same three models (bottom rows) during the historical period 1956-2005. Models shown: MIROC5, GFDL CM3 and MRI. For clarity, the strength of the ENSO-related area is saturated when exceeding the colorscale and its value is indicated at the top of each panel, together with D and ARI from HadISST or ERA40+Interim for each of the model networks	78
33	Sea surface temperature link maps from the ENSO-related area in black for HadISST and the three sample models during the historical period 1956-2005. Models shown: MIROC5, GFDL CM3 and MRI	79
34	Trend anomaly maps for boreal winter in the second half of the 21st century. Anomalies are computed by removing the global mean trend calculated over the months of December to February and indicated in Table 3 from each grid cell. + and • indicate agreement in more than 90% and 70% of models in the sign of the trend anomaly slope. (a) SST averaged across models over 2051-2100. (b) As in (a) but for rainfall. The units are C°/year for SST and (mm/day)/year for precipitation	80
35	Metric D versus ARI for climate model networks during the period 2051-2100. (a) Sea surface temperature. (b) Precipitation. All networks are referenced to the corresponding integration over the historical period. Three levels of noise-to-signal ratios γ are also indicated. D and ARI between HadISST and other sea surface temperature proxies, and ERA40+Interim and other precipitation reanalyses are repeated to provide context	81
36	Sea surface temperature strength maps for two members of the CanESM2 model in the historical period (1956-2005) on top, and in the 21st century (2051-2100) at the bottom. For clarity, the strength of the ENSO-related area is saturated when exceeding the colorscale and indicated in each panel. In the future projections D and ARI from the corresponding historical member are also specified	82
37	Trend anomaly maps for boreal winter in the 22nd and 23rd centuries. Anomalies are computed by removing the global mean trend calculated over the months of December to February and indicated in Table 4 from each grid cell. + and • indicate agreement in more than 90% and 70% of models in the sign of the trend anomaly slope. (a) Sea surface temperature (SST) averaged across models over 2101-2150. (b) Rainfall averaged across models over 2101-2150. (c) As in (a) but for 2151-2200. (d) As in (b) but for 2151-2200. (e) As in (a) but for 2201-2250. (f) As in (b) but for 2201-2250. (g) As in (a) but for 2251-2300. (h) As in (b) but for 2251-2300. The units are C°/year for SST and (mm/day)/year for precipitation	84

38	Metric D versus ARI for seven climate model networks from 2051 to 2300 over five consecutive 50-year periods, from 1 to 5. (a) Sea surface temperature. (b) Precipitation. All networks are referenced to the corresponding integration over the historical period. Three levels of noise-to-signal ratios γ are also indicated	85
39	Sea surface temperature (a-d) and precipitation (e-h) strength maps for two models (left column CCSM4, right column MPI) in the historical period (1956-2005) and in the future (2251-2300). For each variable the first row corresponds to the historical experiments. For clarity, the strength of the ENSO-related area is saturated when exceeding the colorscale and indicated at the top of each panel. D and ARI metrics of the future projections from the corresponding historical member are also included	86
40	Link maps for sea surface temperature (a-b) and precipitation (c-d) from the ENSO-related area in black for two models for which the ENSO projected strength evolves in opposite ways. CCSM4 is shown on the left column and MPI on the right. Maps are calculated over the 2251-2300 period . . .	87
41	Variance of the cumulative anomalies of the ENSO area in DJF in the models and HadISST over 1956-2005 in red, and in the models over 2251-2250 in blue. For HadISST the time series is highly correlated (coefficient 0.94) with the Niño3.4 index defined as the average of SST anomalies from $5^{\circ}S$ to $5^{\circ}N$, and from 120° to $170^{\circ}W$. Error bars around the mean variance over 50 years are determined using a 20-year sliding window, and provide a measure of the decadal modulation of ENSO in the models over the periods considered.	87
42	Maps of area strength of sea surface temperature networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded	92
43	Link maps from the ENSO related area (in black) for sea surface temperature networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown	93
44	Maps of area strength for precipitation networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded	94
45	Link maps from the ENSO related area (in black) for precipitation networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown	95

46	Maps of area strength for the sea surface temperature networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100). For each model, the ensemble member shown projects into the future the historical counterpart in Fig. 42. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded	96
47	Link maps from the ENSO related area (in black) for sea surface temperature networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100) for the ensemble members in Fig. 46	97
48	Maps of area strength for the precipitation networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100). For each model, the ensemble member shown is the projection into the future of the historical counterpart in Fig. 44. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded	98
49	Link maps from the ENSO related area (pictured in black) for the precipitation networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100) for the ensemble members in Fig. 48	99
50	Maps of area strength for the sea surface temperature networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300). For each model, the ensemble member shown projects into the future the historical counterpart in Fig. 42. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded	100
51	Link maps from the ENSO related area (in black) for sea surface temperature networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300) for the ensemble members in Fig. 50	101
52	Maps of area strength for the precipitation networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300). For each model, the ensemble member shown projects into the future the historical counterpart in Fig. 44. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded	102
53	Link maps from the ENSO related area (in black) for the precipitation networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300) for the ensemble members in Fig. 52 .	103
54	Areas identified using three different cell-level networks. α was set to 1×10^{-3} . Data set: HadiSST 1956-2005	106

55	ARI between a reference network constructed using $\alpha = 1 \times 10^{-3}$ and networks constructed using different α values	106
56	Correlogram between two climate time series for a lag range of ± 12 months. We show the significant correlations for a false discovery rate $q = 10^{-3}$ with red. The error bars correspond to \pm one standard deviation, as estimated by Eq. (15).	115
57	(A) The identified domains. The color of each domain corresponds to the connected component it belongs to (the blue and green nodes belong to two different poles of the same component). (B) Color map for domain strength. The strength of ENSO (domain E) is shown at the top. (C) Edges to and from ENSO (shown in black). (D) The climate network. The color of each edge represents the corresponding cross-correlation. (E) The lag range associated with each edge. (F) Examples of lag-consistent triangles. .	121
58	(A),(B) The first two components of EOF analysis. (C) Communities identified by OSLOM. Each community has a unique number and color. (D) Areas identified by spatial clustering.	122
59	Three domain-level network communities for each scan. The first corresponds to the default-mode network, the second to the occipital network, and the third to the motor/somatosensory network.	125
60	The domains of the backbone network for each hemisphere and scan. The color of each domain is randomly assigned (overlaps are shown in black). .	126

SUMMARY

Spatio-temporal data have become increasingly prevalent and important for in many scientific fields (e.g., climate, systems neuroscience, seismology) and enterprises (e.g., geo-tagged tweets). Such data are typically embedded in an arbitrary grid. The grid cells, however, do not correspond to functionally distinct units. One major task is to identify the distinct semi-autonomous functional components of the system and to infer their interconnections.

Common computational analysis methods for such data include standard time series analysis, clustering, community detection, and multivariate statistical methods (e.g., PCA/ICA). However, as we also demonstrate using synthetic data, each of these classes of methods have important limitations in terms of accuracy and flexibility.

In this thesis, we propose two methods that first identify the functional components of a spatio-temporal system as spatially contiguous sets of grid cells, homogeneous to the underlying field. At a second step, an edge inference process identifies the possibly lagged and weighted connections between the system's components, applying a multiple-testing process controlling for the rate of false positives. The inferred network is modeled as a weighted and directed graph. The weight of an edge accounts for the magnitude of the interaction between two components; the direction (and lag) associated with each edge accounts for the temporal ordering of the interactions between the system's components.

The first method, geo-Cluster, infers the spatial components as "areas". An area is a spatially contiguous, non-overlapping, set of grid cells that satisfy a homogeneity constraint in terms of their average pair-wise cross-correlation. However, in real physical systems the underlying physical components might not have crisp boundaries (i.e., they might overlap). To account for this we also propose δ -MAPS, a method that first identifies the epicenters

of activity of the functional components of the system and then creates domains - spatially contiguous, possibly overlapping, sets of grid cells that satisfy the same homogeneity constraint.

The proposed framework is applied in climate science and neuroscience. In the context of climate we show how such methods can be used to infer climate shifts, evaluate cutting edge climate models and identify lagged relationships between different climate regions. In the context of neuroscience, the method is applied to resting state fMRI data and successfully identifies well-known "resting state networks" as well as a few areas that are strongly interconnected to each other, forming the backbone of the functional cortical network.

Chapter I

INTRODUCTION

Many real world systems are modeled as an ensemble of distinct components that are associated via a complex set of connections. In some systems both the elements and their connections are obvious (e.g., Internet routers as nodes, cables between routers as edges). In others, the underlying mechanisms for remote connections are unknown a priori (e.g., social networks) and it is non-trivial to identify the distinct functional components of the system (e.g., functional regions in the human brain). This is usually the case with systems embedded in a spatio-temporal field.

In recent years, spatio-temporal data have become increasingly prevalent and important for in many scientific fields (e.g., climate, systems neuroscience, seismology) and enterprises (e.g., geo-tagged tweets). Such data are typically embedded in an arbitrary grid. The grid cells, however, do not correspond to functionally distinct units. One major task is to identify the distinct semi-autonomous components of the system. A second is to infer the strength of their (potentially lagged) interconnections.

A typical approach to study spatio-temporal systems is to model them as networks. Typically, the grid cells are the nodes of the network and the edges of the network correspond to statistically significant linear [153] or non-linear [53] relationships between the grid cell time series. These networks are modeled either as binary [165] or weighted graphs [68]. Such methods have been successfully employed to forecast El Niño events [105], uncover interesting global-scale patterns responsible for the transfer of energy throughout the oceans [52], investigate changes in the network structure due to neurobiological disorders [138] and many more. The main drawback of such an approach is that the size and number of the nodes (i.e., grid cells) are arbitrarily determined by the measurement technique and

do not correspond to functionally distinct units.

1.1 Dimensionality reduction methods for spatio-temporal data

To uncover the functional components of a spatio-temporal system, it is necessary to identify the dimensionality in the spatial domain. This can be accomplished through the use of spatial dimensionality reduction techniques.

A common approach to reduce the dimensionality of a spatio temporal system is through multivariate statistical methods. Examples of such methods include Principal Component Analysis (PCA) [92], also known as Empirical Orthogonal Function (EOF) analysis [167], and Independent Component Analysis (ICA) [89]. PCA (standard or rotated) aims to decompose the observed data into orthogonal vectors (i.e., the principal components) of high energy content in terms of the variance of the signal. Known drawbacks of PCA include the fact that lower variance components are masked by higher variance ones, and so the analysis is typically limited to the first one-two principal components, as long as they explain most of the variance. Further, the orthogonality between PCA components complicates the interpretation of the results making it difficult to identify the distinct functional components and separate their effects [50]. ICA separates a mixed signal into independent, non-Gaussian components. In contrast to PCA there is no orthogonality constraint imposed on the identified components. However, one cannot determine the variance, sign, or the correct ordering of the independent components. In other words, ICA does not provide a relative significance for each component and the number of independent components should be chosen based on some additional information about the underlying system. Finally, an independent/principal component does not represent a distinct functional component; it is the mixture of many functional components.

Another broad family of spatio-temporal dimensionality reduction methods is based on

clustering [60, 90]. Examples of clustering algorithms include region growing [104], partitioning [139], hierarchical [24], spectral [160] and probabilistic [83] methods. The functionality and scope of each method differs but they share some common characteristics. For instance, every grid cell needs to belong to a cluster while the actual number of clusters is often required as an input parameter. Further, the identified clusters are non-overlapping and might not be spatially contiguous. In particular, the lack of spatial contiguity makes it hard to distinguish between correlations due to spatial diffusion (or dispersion) phenomena from correlations that are due to remote interactions between clusters. Relevant to clustering are community detection techniques [8, 145], which are applied on the cell-level network directly. In contrast to clusters, communities can be overlapping [4, 116], however there is no spatial contiguity constraint. Further, community detection methods do not decouple the identification of the functional components, to the connections that these have with each other. Two components in the same community might have different connectivity patterns to the rest of the network.

1.2 A framework for the analysis of spatio-temporal systems

In this thesis, we propose a framework that first identifies the distinct semi-autonomous components of a spatio-temporal system as spatially contiguous clusters of grid cells. At a second step, the (possibly lagged) interactions between them are inferred and their magnitude is assessed.

1.2.1 *geo-Cluster*

In detail, we first propose *geo-Cluster*, a method that first infers the spatial components of the underlying system as “areas”. An area is a spatially contiguous, non-overlapping, set of (two or more) grid cells that satisfy a homogeneity constraint based on their average pairwise cross-correlation. For parsimony reasons the proposed method aims to maximize the size of the identified areas. The method requires a single parameter which determines the minimum degree of homogeneity of the grid cells in each area. Next, a complete weighted

network between the identified areas is inferred, modeling the functional relationships between them. The weight of an edge corresponds to the covariance between the area time series, accounting for the power of the signal of each area as well as the correlation between the area time series.

The proposed method has been shown to be robust to noise, the resolution of the underlying grid, the parameter that determines the minimum degree of homogeneity in an area, and the metric used to quantify the similarity between the grid cell time series. *geo-Cluster* has been extensively applied to climate data to investigate climate shifts and to construct interdependent networks [71] between different climate domains. Further, the method is applied to evaluate cutting edge climate models assessing their ability to reproduce the climate in the past and investigating the model trajectories under a future climate warming scenario.

1.2.2 δ -MAPS

In real physical systems the underlying spatial components might not have crisp boundaries [63] and their interactions might not be instantaneous. To this end, we propose δ -MAPS; a method that identifies spatially contiguous and possibly overlapping components referred to as “domains”, and identifies the lagged functional relationships between them. Informally, a domain is a spatially contiguous region that somehow participates in the same dynamic effect or function. The latter will result in highly correlated temporal activity between grid cells of the same domain. Thus, δ -MAPS first identifies the epicenters of activity of a domain. Next, it identifies a domain as the maximum possible set of spatially contiguous grid cells that include the detected epicenters and satisfy a homogeneity constraint (based again on the average pair-wise correlation of the grid cells in the domain’s scope). After identifying the domains, δ -MAPS infers a functional network between them. The proposed network inference method examines the statistical significance of each lagged correlation between two domains, applies a multiple-testing process to control the rate of

false positives, infers a range of potential lag values for each edge, and assigns a weight to each edge reflecting the magnitude of interaction between two domains. δ -MAPS is related to clustering, multivariate statistical techniques and network community detection. However, as we discuss and also show with synthetic data, it is also significantly different, avoiding many of the known limitations of these methods.

We illustrate the application of δ -MAPS on data from two domains: climate science and neuroscience. First, the sea-surface temperature (SST) climate network identifies some well-known *teleconnections* (such as the lagged connection between the El Niño Southern Oscillation and the Indian Ocean). Second, the analysis of resting state fMRI cortical data confirms the presence of known functional resting state networks (default mode, occipital, motor/somatosensory and auditory), and shows that the cortical network includes a backbone of relatively few regions that are densely interconnected.

1.3 Organization

The thesis is organized as follows. In Chapter II we present related work on network inference and dimensionality reduction techniques for spatio-temporal data. Using a synthetic data set, in which the functional components and their interconnections are known, we contrast such methods to the proposed framework and identify key differences and limitations for each method. In Chapter III we propose *geo-Cluster* [66], provide robustness results and show example applications in the field of climate science. In Chapter IV we provide an extensive application of *geo-Cluster* to evaluate cutting edge climate models [67]. In Chapter V, we propose δ -MAPS [65]. We compare δ -MAPS to the most common dimensionality reduction methods and show its application on the fields of climate and neuroscience. Finally, Chapter VI provides the main conclusions from this thesis and an outlook for future work.

Chapter II

RELATED WORK

There exist several methods to analyze spatio-temporal data, from simple time series analysis, to dimensionality reduction methods (e.g., PCA/clustering), to network based methods. The proposed framework combines ideas from the latter two approaches. In the following, we first introduce a synthetic data set in which both the functional components of the system and their interactions are a priori known. We shall use this data set to illustrate limitations of existing dimensionality reduction approaches. Next, we provide a brief overview of the network based approach. We conclude by presenting various dimensionality reduction techniques, contrasting them to the proposed approach, in terms of their ability to uncover the functional components of a spatio-temporal system and their interactions.

2.1 Synthetic Data Generation

To better highlight the limitations and differences of various dimensionality reduction methods, we first introduce an example in which we know both the dimensionality of the spatio-temporal system as well as the interactions between the different components of the system.

We construct five domains (modeled as spatially contiguous regions) on a 50×70 spatial grid. Each domain i is associated with a “mother” time series $y_i(t)$, ($i=1 \dots 5$). To make the experiment more realistic in terms of autocorrelation structure and marginal distribution, each $y_i(t)$ is a real fMRI time series with length $T=1200$ (see Section 5.6). The five mother time series $y_i(t)$ are uncorrelated (absolute cross-correlation < 0.05 at all lags), and they are normalized to zero-mean, unit-variance. To create correlations between domains (i.e., domain-level edges), we construct five new time series $x_i(t)$ based on linear combinations of two or more mother time series. For instance, if we set $x_i(t) = (1 - \alpha)y_i(t) + \alpha y_j(t + \tau)$

with $0 < \alpha < 1$ and $x_j(t) = y_j(t)$, domains i and j become positively correlated at a lag τ ; the correlation increases with α . The time series x_i are again normalized to zero-mean, unit-variance. We then scale the time series of domain i by a factor $\sqrt{s_i}$ to control the variance of each domain ($\text{Var}[x_i(t)] = s_i$).

For simplicity, each domain is a circle with radius r_p . A domain has a “core region” with the same center and radius $r_c < r_p$; the core is supposed to be the epicenter of that domain. Every point in the core has the same signal $x_i(t)$ (before we add random noise). Outside the core, the signal attenuates at a distance d from the center of the domain as follows:

$$x_i(t) = \sqrt{f(d)} x_i(t), f(d) = \frac{r_p - d}{r_p - r_c}, r_c \leq d \leq r_p. \quad (1)$$

The parameters of the five synthetic domains are shown in Table 1. The domains differ in terms of size and power (variance). The spatial extent of the domains is shown in Fig.1-A; domains 1 and 3 overlap with domain 2, while domains 4 and 5 also overlap to a smaller extent. Further, there is a strong and lagged anti-correlation between domains 1 and 3, a weaker positive correlation at zero-lag between domains 4 and 5, and an ever weaker positive correlation at zero-lag between domains 3 and 5. The edges of the domain-level network are also shown in Fig.1-A.

Table 1: Synthetic domain generation parameters.

ID	r_c	r_p	s_i	$x_i(t)$
1	2	10	16	$x_1(t) = 2/3y_1(t) - 1/3y_3(t + 15)$
2	4	14	11	$x_2(t) = y_2(t)$
3	2	10	16	$x_3(t) = y_3(t)$
4	0.5	5	9	$x_4(t) = 3/4y_4(t) + 1/4y_5(t)$
5	1	7	6	$x_5(t) = 4/5y_5(t) + 1/5y_3(t)$

2.2 Network Based Methods

Spatio-temporal data are usually embedded in a two or three dimensional grid; each grid cell contains time series of measurements for a given variable. Such a data set can be naturally modeled as a network. The grid cells play the role of the nodes in the network. The

edges of the network are inferred based on statistically significant relationships between the grid cell time series [52, 106, 117]. The network can be modeled either as a binary [153] or weighted [68] graph.

All these methods require a statistical test to distinguish between significant and non-significant edges. Naive approaches to the problem include using a fixed threshold [153] or requiring the network to have a fixed density [142]. More sophisticated approaches such as using surrogate time series also exist (see e.g., [100, 107, 127]). However, methods based on surrogate time series are computationally expensive (compared to the naive approach) and might not scale for finer-scale resolution data.

The network based approach has been successful in many fields. For example, complex network approaches have been used to forecast El Niño events [105], map brain regions that are most likely to be affected by pathological changes [45], identify structures responsible for the transfer of energy in oceans [52], show how some diseases (such as Alzheimer's) can affect the functional structure of the brain network [147] and many more. The main drawback of such an approach is that the size (and number) of the nodes are arbitrarily determined by the measurement technique and do not correspond to functionally distinct units. To counter this problem there has been an effort to identify *modular* networks. From the network perspective, a module (or *community*) corresponds to a set of grid cells highly interconnected to each other and less connected to the rest of the world. The identification of communities in networks is a concept similar to clustering and will be discussed further in Section 2.3.4.

2.3 Dimensionality Reduction methods

In this section we provide an overview of dimensionality reduction methods used to analyze spatio-temporal data. For each family of methods we show their limitations, when the objective is to identify the functional components and their interconnections in our synthetic data set.

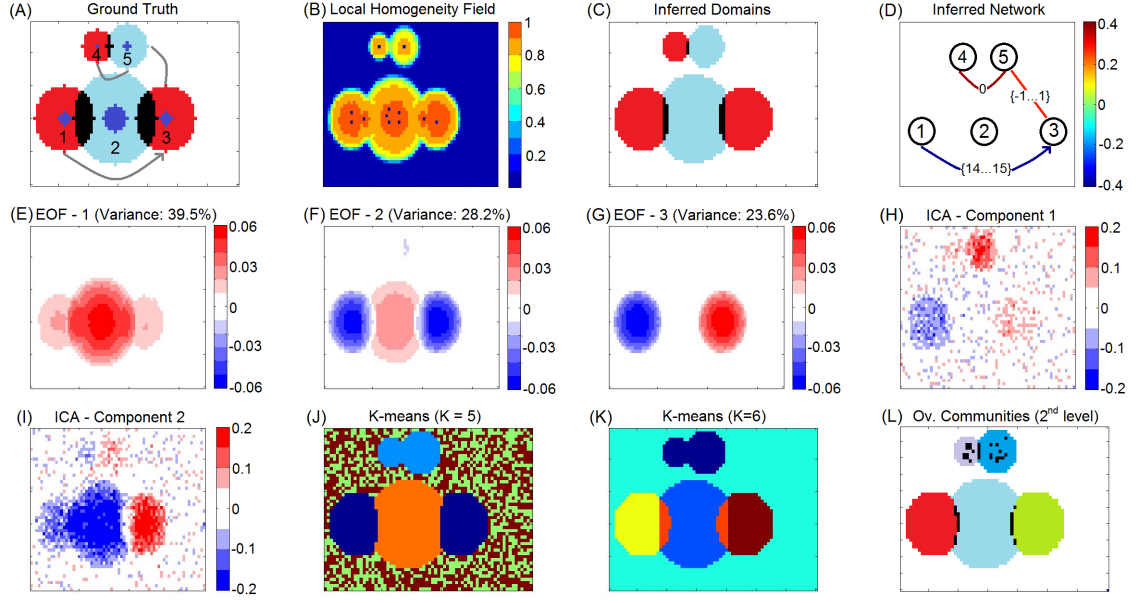


Figure 1: **A:** The five ground-truth domains. Adjacent domains have different colors, overlapping regions shown in black, and the core of each domain is in blue. The three constructed edges are shown in gray lines. **B:** The homogeneity field $\hat{r}_K(i)$ at each cell. The identified seeds are shown in blue. **C:** The inferred domains: adjacent domains have different colors and overlaps are shown in black. **D:** The inferred domain-level network: the color map refers to the edge correlation. The lag associated with each edge is also shown. **E,F,G:** The first three EOF (PCA) components. The variance explained by each component is shown at the top of each figure. **H,I:** The two ICA components. **J,K:** K-means clustering. **L:** The second hierarchical level of community structure as identified by OSLOM: each community has a distinct color and overlaps are shown in black.

2.3.1 Principal Component Analysis

Principal Component Analysis (PCA), also known as Empirical Orthogonal Function (EOF) analysis, is one of the oldest techniques used to analyze spatio-temporal data. PCA aims to decompose the original set of variables into a new set of (principal) components that capture most of the observed variance of the data through a linear combination of the original variables. The identified components are orthogonal to each other and each component is assigned a value equal to the total variance explained by it.

PCA assumes that the dominant patterns are orthogonal in space and time (which is not necessarily true, see [132] for a case relevant to climate). To overcome this problem alternative methods (e.g., rotated PCA) exist [164] but require more user defined parameters and sometimes split a single pattern into two different ones [167]. An interesting analysis on how PCA results can be misleading is presented in [50].

We apply PCA using Matlab's PCA toolbox. Fig. 1-E,F,G show the first three principal components, which collectively account for about 90% of the total variance. A first observation is that domains 4 and 5 are not even visible in these components – they only appear in the next two components, which account for about 5% of the variance each. This is because domains 4 and 5 are smaller and have lower variance. This is a general limitation of PCA: the variance of the analyzed field can be dominated by a small number of “modes of variability”, completely masking smaller/weaker regions of interest and their connections. Second, the first three components do not provide a consistent evidence that domains 1 and 3 are strongly anti-correlated; this is due to their lagged correlation, which is missed by PCA. Third, the first component, which accounts for 40% of the total variance, can be misinterpreted to imply that domain 2 is somehow positively correlated with domains 1 and 3, even though it is actually generated by an uncorrelated signal. This is due to the overlap of domain 2 with domains 1 and 3.

2.3.2 Independent Component Analysis

A method typically used by the neuroscience community is Independent Component Analysis (ICA) [89]. ICA defines a model for the observed data; in the model the data are assumed to be linear mixtures of some unknown latent variables (the independent components). The mixing system is also unknown. The unknown latent variables are assumed to be non-Gaussian and mutually independent. In general, given the observed signals x , the goal is to identify the mixing matrix A and the independent components s such that $x = As$. In contrast to PCA there is no orthogonality constraint on the independent components but there is no way to determine the number, the variance, the sign, or the correct ordering of the independent components. Thus, one should rely on empirical knowledge to identify independent components that correspond to specific functional modules.

We apply ICA on the synthetic data using Matlab’s FastICA toolbox. To help ICA perform better, we actually specified the right number of independent components, which is two (domains 1,3,4,5 are indirectly correlated – domain 2 is not correlated with any other). The two independent components are shown in Fig. 1-H,I. Note that only a rough “shadow” of each domain is visible. Domains 1 and 3 appear in different colors, providing a hint that they are anti-correlated, while domains 3 and 5 appear in the same color because they are positively correlated. Overall, however, the components are quite noisy and it would be hard in practice to discover the functional structure of the underlying system if we did not know the ground-truth. The results are even harder to interpret when we request a larger number of components.

2.3.3 Clustering based methods

A broad family of dimensionality methods is based on clustering [13, 139, 152, 160]. None of these methods though guarantee that the identified clusters are spatially contiguous. As we show next, spatial contiguity is an attractive property since it will enable us to differentiate network nodes from large scale networks of nodes [17].

We apply the most well-known clustering method, *k-means*, on our synthetic data. As commonly done with correlation-based clustering, the distance between two cells i and j is determined by the maximum absolute correlation across all considered lags, as $1 - |r_{i,j}^*|^1$. Fig. 1-J,K shows the resulting clusters for $k=5$ (the number of synthetic domains) and 6, respectively. For $k=5$, domains 1 and 3 form a single cluster because of their strong anti-correlation; the same happens with domains 4 and 5. Further, two of the five clusters (green and brown) cover just noise. The situation changes completely when we request $k=6$ clusters. In that case, the overlapping regions in domain 2 form a single cluster, while domains 1 and 3 are separated in different clusters.

More similar to the proposed framework is the notion of spatially contiguous clustering. Identifying spatially contiguous clusters is a problem that arises in many fields, from image processing [69], to geographical sciences [55], to studying the Earth's climate [66] and the human brain [44]. In general there exist two approaches to find spatial clusters. The first approach is a semi-supervised approach where after the initial clusters are identified subsequent (supervised) changes are made to merge them into spatially contiguous regions [55]. Our focus here are unsupervised approaches, where the spatial contiguity criterion is incorporated into the clustering algorithm.

A typical approach to identify spatially contiguous clusters is agglomerative hierarchical clustering. In such an approach each grid cell forms its own cluster and clusters are iteratively joined, according to some distance measure, if they are spatially adjacent. Many distance measures have been proposed in the literature. For example, in [79] the authors evaluate a family of three different distance measures (single, average and complete linkage) to cluster U.S. presidential election data into different regions. Essentially this

¹In detail, the Pearson correlation between grid cells i, j and lag τ is given by $r_{i,j}(\tau) = \frac{\sum_{t=1}^{T-\tau} (x_i(t) - \tilde{\mu}_i)(x_j(t+\tau) - \tilde{\mu}_j)}{T\tilde{\sigma}_i\tilde{\sigma}_j}$, T being the time series length and $\tilde{\mu}_i, \tilde{\sigma}_i$ their empirical mean and standard deviation respectively. $|r_{i,j}^*| = \arg \max_{\tau \in \{-\tau_{max} \dots \tau_{max}\}} |r_{i,j}(\tau)|$, with τ_{max} being the maximum lag.

hierarchical clustering process constructs a dendrogram; by “cutting” horizontally the dendrogram at a level of our choice we obtain the resulting clusters. In the case of the human brain, hierarchical clustering has been applied to provide a full brain parcellation [111] and to test clusters (rather than voxels) for activation during task based experiments [81].

Another popular clustering method with applications in fMRI is NCUT [44, 131]. NCUT is a graph based clustering method; edges are removed iteratively until a pre-specified number of clusters is reached. The edges are removed such as to maximize the similarity between the elements in the same cluster while maximizing the dissimilarity between elements in different clusters. Finding the optimal edge to remove is NP-Complete and the method relies on an approximate solution (found by solving a generalized eigenvalue problem). One of the drawbacks of NCUT is that it is biased to identify clusters of similar size (for a detailed description of the limitations of spectral clustering methods we refer the reader to [112]). In [150] the authors compare NCUT to an agglomerative hierarchical clustering that uses Ward’s distance [169]. They show that the latter performs better both in terms of reproducibility (i.e., sensitivity to noise) as well as in terms of accuracy.

Another group of clustering methods is based on the concept of region growing [2]. Typically, region growing methods start with a number of pre-specified seed regions (a seed region contains only one grid cell). These regions grow by including grid cells similar to them, until a homogeneity criterion is reached. Similar to our approach region growing methods are based on the intuition that neighboring grid cells should have similar values. In contrast to the proposed method, there is no merging of regions while these are growing. Selecting the location of seed regions will affect the outcome of the clustering algorithm and a couple of different approaches exist. For example, in [104] the authors propose that all grid cells form a seed region. Having identified a seed region for each cell in the grid they iteratively remove (and keep) the largest regions up to the point that only regions of size less than an arbitrary threshold remain. In [24] the authors use the concept of stability maps to select the seed regions and at a second step they use a hierarchical agglomerative

clustering to identify functional regions in the cortex.

In contrast to *geo-Cluster*, all these clustering methods need as an input the number of clusters to identify. Moreover, such clustering methods identify spatially contiguous regions even if the underlying field is composed of noise. Further, many of these methods (e.g. NCUT) can be applied only when the distances between grid cells are positive. If the distance between the grid cell time series is captured by a measure which also takes negative values (e.g., Pearson correlation with is the norm) then the similarity matrix has to be thresholded to remove them. Similarly to *geo-Cluster*, the borders of the clusters are crisp and no overlaps are allowed.

An alternative to clustering are edge detection or border detection methods. Border detection techniques are based on the idea that “pixels” or grid cells representing the same object should have a similar value yet distinct from pixels belonging to another area. In [15] for example, the authors use a graph based border detection technique to extract homogeneous regions from raster data. In [38, 76, 170] the idea of border detection is applied to fMRI data where the authors try to delineate the borders of functional areas. Border detection techniques are known to be sensitive to localized patches of noise in the data.

2.3.4 Community Detection Methods

All of these clustering methods suppose that the identified clusters are independent in the spatial domain with their boundaries well defined. In reality, the borders of the clusters might not be clearly demarcated. Some grid cells belonging to one cluster can belong to other clusters as well. For example, when we study the Earth’s climate we are interested to identify functional domains (e.g. the El Niño Southern Oscillation). Such domains have identifiable effects in the temperature anomaly field. One could not claim that strict borders exist in the gradients of the temperature as to allow the definition of “crisp” clusters. In the field of neuroscience there is further evidence of the existence of overlapping clusters.

Quoting Fornito et al.: “*A model of neuronal architecture that allows for overlapping modules offers a more realistic brain-network organization (for instance, cortical association areas are known to have a role in multiple networks)*” [63]. Further, other studies suggest that cognitive functions are organized into segregated and overlapping networks [54, 82].

To our knowledge, the only method that allows for overlapping partitions is community detection. A community is a set of nodes that are highly interconnected to each other, while having fewer connections to the rest of the network. There are numerous methods to identify communities, for a review we refer the reader to [64]. There are also many applications of community detection in the fields of climate science and neuroscience. For example in [145] the authors use community detection techniques to evaluate climate models while in [143] the authors propose the use of communities as informative predictors in lieu of climate indices. In [118] the authors apply a wide variety of community detection methods in fMRI data. At a second step, using a map of task-based activations, they map these communities to specific cognitive functions. Many authors [21, 176] have also suggested that the community structure of the human brain deteriorates (i.e. becomes less modular) as a person gets older. For a comprehensive review of applications of community detection methods in neuroscience we refer the reader to [137].

Overlapping communities are a “natural” extension to the classic definition of a community. The main premise is that an individual can belong to more than one communities (e.g. work, family etc.). There exist several approaches to identify overlapping communities (e.g. [4, 59, 116]). One of the first applications of overlapping community detection in spatio-temporal data (and more specifically in resting state fMRI data) can be found in [175]. The authors test the capability of the proposed methodology to uncover overlapping communities in the resting state network. In [171] the authors investigate the overlapping community structure in the structural brain network and show that the identified communities can be mapped to well-known brain systems. To the best of our knowledge none

of these methods guarantees that the identified communities form spatially contiguous regions.

We apply a state-of-the-art overlapping community detection method, referred to as OSLOM [103], with the default parameter values. The input to OSLOM is a positively weighted graph: each vertex is a grid cell and an edge between vertices i and j corresponds to the maximum absolute cross-correlation $|r_{i,j}^*|$ across all lags of interest. Absolute correlations less than 30% are considered insignificant and the corresponding edges are pruned.² As most community detection methods, OSLOM does not distinguish between positive and negative correlations. OSLOM provides a hierarchy of communities. When applied to our synthetic data, the first level of hierarchy (not shown) simply groups together domains 1,2,3 in one community (even though domain 2 is uncorrelated with domains 1 and 3), and domains 4,5 in another community. The connection between domains 3 and 5 is missed. The second level of hierarchy is shown in Fig. 1-L. Overall, OSLOM does a better job than PCA/ICA/clustering in detecting the spatial extent of each domain. A small overlap between domains (1,2) and (2,3) is discovered but to a smaller extent than δ -MAPS (the results of δ -MAPS are discussed in more detail in Section 5.4). However, a community in OSLOM is not constrained to be spatially contiguous. This is the reason we see some black dots in regions 4 and 5; these are non-contiguous overlaps between the communities that correspond to these two domains. Thus, a community may group together two regions that are, first, not spatially contiguous, and second, different in terms of how they are connected to other regions.

²We have experimented with other pruning thresholds between 20%-50% and the results are very similar at the first two hierarchy levels.

Chapter III

GEO-CLUSTER: SPATIO-TEMPORAL NETWORK ANALYSIS **FOR STUDYING CLIMATE PATTERNS**

3.1 Introduction

Network analysis refers to a set of metrics, modeling tools and algorithms commonly used in the study of complex systems. It merges ideas from graph theory, statistical physics, sociology and computer science, and its main premise is that the underlying topology or network structure of a system has a strong impact on its dynamics and evolution [114]. As such it constitutes a powerful tool to investigate local and non-local statistical interactions.

The progress made in this field has led to its broad application; many real world systems are modeled as an ensemble of distinct elements that are associated via a complex set of connections. In some systems, referred to as structural networks, the underlying network structure is obvious (e.g. Internet routers as nodes, cables between routers as edges). In others, the underlying mechanisms for remote connections between different subsystems are unknown *a priori* (e.g. social networks, or the climate system); still, their effects can be mapped into a functional network. An extensive bibliography for applications of network analysis can be found in [113].

By quantifying statistical interactions, network analysis provides a powerful framework to validate climate models and investigate teleconnections, assessing their strength, range, and impact on the climate system. The intention is to uncover relations in the climate system that are not (or not fully) captured by more traditional methodologies used in climate science [49, 43, 1, 73, 72, 62, 10, 11], and to explain known climate phenomena in terms of the underlying network's structure and metrics.

Introductions to the application of network analysis in climate science are presented in

[142] and [156]. We can classify the prior work in this area in three distinct approaches. A first approach assigns known climate indices as the nodes of the network [154, 148, 168]. By studying the collective behavior of these nodes, it has been possible to investigate their relative role over time and to interpret climate shifts in terms of changes in their relative strength. This approach is obviously sensitive to the initial selection of network nodes, and it cannot be used to discover new climate phenomena involving other regions.

A second, and more common, approach represents the nodes of the climate network by grid cells in the given climate field. Specifically, each grid cell is represented by a node, and edges between nodes correspond to statistically significant relations based on linear or nonlinear correlation metrics [153, 53]. In this approach, it is common to prune edges whose statistical significance is below a certain threshold, and to assume that all remaining edges are equally “strong”, resulting in an unweighted network [157, 53, 142]. This approach has been used to study teleconnections, uncover interesting global-scale patterns responsible for the transfer of energy throughout the oceans, and analyze relations between different variables in the atmosphere [157, 155, 173, 52, 51]. A limitation of this approach is that it results in a very large number of network nodes (all cells in a grid map), and these nodes cannot be used to describe parsimoniously any identified climate phenomena.

The third approach focuses on the community structure of the underlying network [115]. A community is a collection of nodes that are highly interconnected, while having much fewer interactions with the rest of the network. Communities can serve as informative predictors in lieu of climate indices [158, 143, 117], while their evolution and stability has also received some attention [142, 144]. Clustering techniques have also been proposed to discover significant geographical regions in a given climate field (again, in lieu of climate indices) [139], and to identify dipoles (i.e., two regions whose anomalies are anti-correlated) and to evaluate their significance [95, 94]. These community-based or clustering techniques, however, do not infer a network of teleconnections between different

communities (clusters), and they do not quantify the intensity of teleconnections between geographically separated regions within the same community (cluster).

In this work, we propose a new method to apply network analysis to climate science. We first apply a novel network-based clustering method to group the initial set of grid cells in “areas”, i.e., in geographical regions that are highly homogeneous in terms of the underlying climate variable. These areas represent the nodes of the inferred network. Links between areas (i.e., the edges of the network) represent non-local dependencies between different regions over a certain time period. These inter-area links are weighted, and their magnitude depends on both the cumulative anomaly of each area and the cross-correlation between the two cumulative anomalies. The similarity of our method to previous community/clustering techniques is that nodes are endogenously determined during the data analysis process. The main differences are that each node corresponds to a distinct geographical region, and these nodes form a weighted network based on the connection intensity that is inferred for each pair of nodes. In other words, the proposed method decouples the identification of the geographical boundary of each network node from the estimation of the connection intensity between different regions.

The proposed method requires a single parameter τ , which determines the minimum degree of homogeneity between cells of the same area. The method is robust to additive noise, changes in the resolution of the given data set, the selection of the correlation metric, and variations in τ . The resulting climate network can be applied, regionally or globally, to identify and quantify relationships between climate areas (or teleconnections) and their representation in models, and to investigate climate variability and shifts. Finally, the proposed method can be extended to investigate interactions between different climate variables.

The rest of this chapter is organized as follows: In Section 3.2 we introduce the data sets analyzed in this work. We describe the climate network construction algorithm and the network analysis metrics in Sections 3.3 and 3.4, respectively. The robustness of the climate

network inference process is examined in Section 3.5. Applications of the proposed method to a suite of reanalyses and model data sets are presented in Section 3.6. A discussion of the main outcome of this work concludes the chapter.

3.2 Data sets

In this section we briefly describe the data sets that are used in the rest of this chapter. For sea surface temperatures (SSTs), we construct and compare networks based on the HadISST [121], the ERSST-V3 [134] and the NCEP/NCAR [93] reanalyses. For precipitation, we rely on CMAP merged data [172] and ERA-Interim reanalysis [46]. We also analyze the SST fields generated by two coupled general circulation models chosen from the CMIP5 archive: the NASA GISS-E2H [80] and the Hadley Center HadCM3 [75]. We select randomly two runs of each model from the “historical run” ensembles [149].

Because the quality of the measurements contributing to the SST reanalyses deteriorates as we move to higher latitudes, we only consider the latitudinal range of $[60^{\circ}N; 60^{\circ}S]$, avoiding sea-ice covered regions. Also, we mostly focus on the period 1979-2005; in the case of HadISST reanalysis, we contrast with the network characteristics during the 1950-1976 interval. Due to space constraints, results are only shown for the boreal winter season (December to February, DJF). When not specified otherwise, all SST data are interpolated (using bilinear interpolation) to the minimum common spatial resolution across all data sets ($2^{\circ} \times 2.5^{\circ}$); for precipitation the resolution is $2.5^{\circ} \times 2.5^{\circ}$.

All climate networks are constructed from detrended anomalies derived from monthly averages of the corresponding climate field. The detrending is done using linear regression and the anomalies are computed after removing the annual cycle.

3.3 Climate network construction

The network construction process consists of three steps. First, we compute the “cell-level network” from the detrended anomaly time series of each cell in the spatial grid. Second,

we apply a novel *area identification algorithm* on the cell-level network to identify the nodes of the final “area-level network”; an area here represents a geographic region that is highly homogeneous in terms of the given climate field. Third, we compute the weight of the edges between areas, roughly corresponding to teleconnections, based on the covariance of the cumulative anomalies of the two corresponding areas. The following network construction method requires a single parameter, τ , which determines the minimum degree of homogeneity between cells of the same area. In the following we describe each step in more detail.

3.3.1 Cell-level network

Consider a climate field $\mathbf{x}(t)$ defined on a finite number of cells in a given spatial grid. The i ’th vector of the climate field is a time series $x_i(t)$ of detrended anomalies in cell i . The length of each time series is denoted by T . We first compute Pearson’s cross-correlation $r(x_i, x_j)$ ¹ between the time series $x_i(t)$ and $x_j(t)$ for every pair of cells i and j . We calculate the correlations at zero-lag, assuming that the physical processes linking different cells result from atmospheric wave dynamics and are fast compared to the *one-month averaging time scale* of the input time series. Considering time-lagged correlations is beyond the scope of this chapter. Instead of using Pearson’s correlation, other correlation metrics could be adopted; in Section 3.5.4 we examine the differences in the resulting network using a rank-based correlation metric.

Most of prior work on climate network analysis applies a cutoff threshold on the correlations $r(x_i, x_j)$ to prune insignificant ones and construct a binary (i.e., unweighted) network between cells; for a recent review see [141]. Fig. 2 shows correlation distributions for four SST reanalysis data sets; note that there is no natural cutoff point to separate significant correlations from noise. We have experimented with methods that first prune insignificant correlations and then construct unweighted networks, and observed that the final area-level

¹Unless specified otherwise, the term “correlation” will be used to denote Pearson’s cross-correlation metric between two time series.

network is sensitive to the significance level at which correlations are pruned. Such sensitivity complicates any attempt to make quantitative comparisons between networks constructed from different data sets (for example networks from observations versus models).

For this reason, in the following we present a method that considers *all* pair-wise cell correlations, without any pruning. Thus, the cell-level network is a *complete and weighted* graph, meaning that every pair of cells is connected but with weighted edges between -1 and 1. This cell-level network is the input to the area identification algorithm, described next.

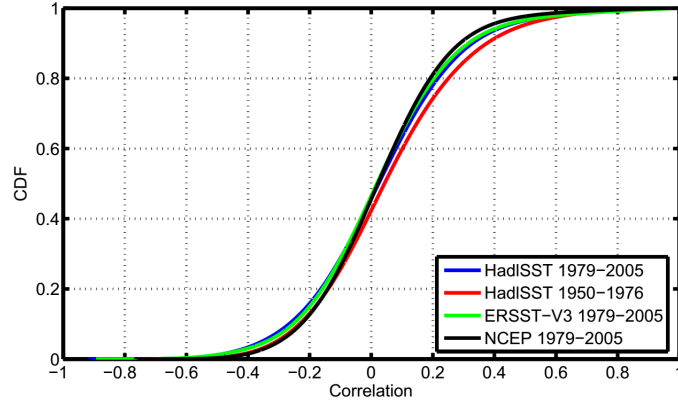


Figure 2: Empirical Cumulative Distribution Functions (CDF) of correlations for the HadISST reanalysis during the 1950-1976 and 1979-2005 periods, and for ERSST-V3 and NCEP data during the 1979-2005 period

3.3.2 Identification of climate areas

A central concept in the proposed method is that of a *climate area*, or simply *area*. Informally, an area A represents a geographic region that is highly homogeneous in terms of the climate field $\mathbf{x}(t)$.

In more detail, we define as *neighbors* of a grid cell i the four adjacent cells of i , and as *path* a sequence of cells such that each pair of successive cells are neighbors. An area A is a set of cells satisfying three conditions:

1. A includes at least two cells.

-
2. The cells in A form a connected geographic region, i.e., there is a path within A connecting each cell of A to every other cell of that area.
 3. The average correlation between all cells in A is greater than a given threshold τ ,

$$\frac{\sum_{i \neq j \in A} r(x_i, x_j)}{|A| \times (|A| - 1)} > \tau \quad (2)$$

where $|A|$ denotes the number of cells in area A .

The parameter τ determines the minimum degree of homogeneity that is required within an area. A heuristic for the selection of τ is presented in Section 3.8; we use that heuristic in the rest of this chapter.

For the climate network to convey information in the most parsimonious way, *the number of identified climate areas should be minimized*. To this end, an area is defined as a maximum cardinality set of cells, that are spatially contiguous, and whose average pairwise correlation is larger than the threshold τ . In Sec. 5.8 we show that this computational problem is NP-Complete, meaning that there exists no efficient way to solve it in practice. Consequently, we have designed an algorithm that aims to *minimize the number of areas* heuristically, based on a so called “greedy” approach [41]. The algorithm consists of two parts. First, it identifies a set of areas; secondly it merges some of those areas together as long as they satisfy the previous three area constraints. A pseudocode describing the algorithm is given in Section 3.9, while the actual software is available at <http://www.cc.gatech.edu/~dovrolis/ClimateNets/>. An example of the area identification process applied to a synthetic grid is illustrated in Fig. 3.

The identification part of the algorithm produces areas that are geographically connected by always expanding an area through neighboring cells. Additionally, the algorithm attempts to identify the largest (in terms of number of cells) area in each iteration by selecting, in every expansion step, the neighboring cell that has the highest average correlation with existing cells in that area. The expectation is that this greedy approach allows the area to expand to as many cells as possible, subject to the constraint that the average correlation

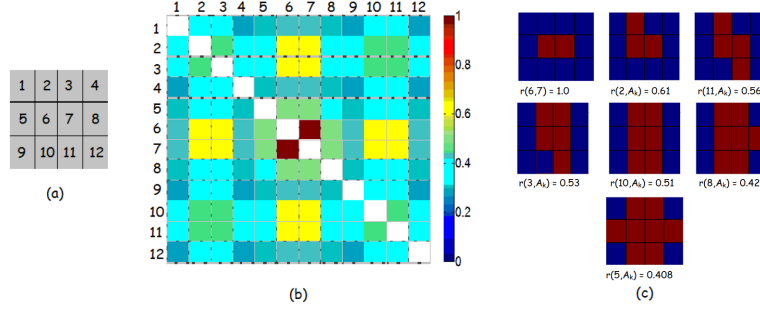


Figure 3: An example of the area identification algorithm. (a) 12-cell synthetic grid. (b) The correlation matrix between cells (given as input). (c) The area expansion process for a given $\tau=0.4$. Cells shown in red are selected to join the area (denoted by A_k). Cells 1, 4, 9 and 12 will not join A_k since they do not satisfy the τ constraint in Eq.2

in the area should be more than τ . It is easy to show that an identified area satisfies the condition given by Eq.2.

Within the set of areas V identified by the first part of the algorithm, it is possible to find some areas that can be merged further, and still satisfy the previous three constraints. Specifically, we say that two areas A_i and A_j can be merged into a new area $A_k = A_i \cup A_j$ if A_i and A_j have at least one pair of geographically adjacent cells and the average correlation of cells in A_k is greater than τ . The second part of the algorithm, therefore, attempts to merge as many areas as possible (see Section 3.9).

Fig. 4 shows the identified areas before merging (i.e., after Part-1 in Section 3.9) and after merging (i.e., after Part-2 in Section 3.9) for the HadISST reanalysis. Fig 4c shows the distribution of area sizes (in number of cells) before and after merging. Area merging decreases substantially the number of small areas (the percentage of areas with less than 10 cells in this example drops from 46% to 10%).

The identified areas represent the nodes of the inferred climate network. We refer to this network as “area-level network” to distinguish it from the underlying cell-level network.

3.3.3 Links between areas

Links (or edges) between areas identify non-local relations and can be considered a proxy for climate teleconnections. To quantify the weight of these links, we first compute for

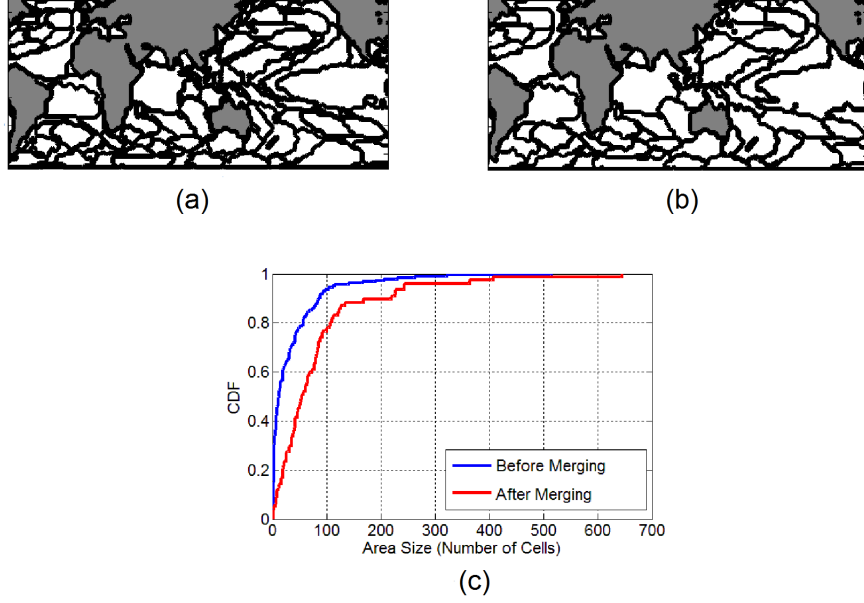


Figure 4: Identified areas in the HadISST 1979-2005 data set ($\tau=0.496$). (a) The 176 areas identified by Part-1 of the area identification algorithm. (b) The 74 “merged” areas after the execution of Part-2. (c) The CDF of area sizes (in number of cells) before and after the merging process

each area A_k the *cumulative anomaly* $X_k(t)$ of the cells in that area,

$$X_k(t) = \sum_{i \in A_k} x_i(t) \cos(\phi_i) . \quad (3)$$

The anomaly time series of a cell i is weighted by the cosine of the cell’s latitude (ϕ_i), to account for the cell’s relative size. As a sum of zero-mean processes, a cumulative anomaly is also zero-mean.

Fig. 5 quantifies the relation between the size of the areas ($\sum_{i \in A_k} \cos(\phi_i)$) identified earlier in the HadISST data set and the standard deviation of their cumulative anomaly. Note that the relation is almost linear, at least excluding the largest 3-4 areas. Exact linearity would be expected if all cells had the same size, their anomalies had the same variance, and every pair of cells in the same area had the same correlation. Even though these conditions are not true in practice, it is interesting that the standard deviation of an area’s cumulative anomaly is roughly proportional to its size.²

²When comparing data sets with different spatial resolution, the anomaly of a cell should be normalized

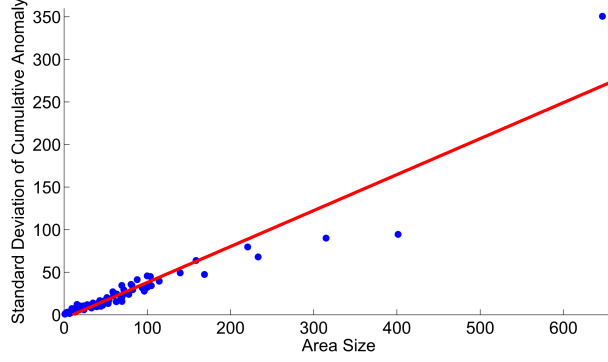


Figure 5: The relation between area size and standard deviation of the area's cumulative anomaly ($R^2 = 0.88$) for the HadISST reanalysis during the 1979-2005 period; $\tau=0.496$

The strength, or weight, of the link between two areas A_i and A_j is captured by the *covariance* of the corresponding cumulative anomalies $X_i(t)$ and $X_j(t)$. Specifically, every pair of areas A_i and A_j in the constructed network is connected with a link of weight $w(A_i, A_j)$,

$$w(A_i, A_j) \triangleq w(X_i, X_j) = \text{cov}(X_i, X_j) = s(X_i) s(X_j) r(X_i, X_j) \quad (4)$$

where $s(X_i)$ is the standard deviation of the cumulative anomaly $X_i(t)$, while $\text{cov}(X_i, X_j)$ and $r(X_i, X_j)$ are the covariance and correlation, respectively, of the cumulative anomalies $X_i(t)$ and $X_j(t)$ that correspond to areas A_i and A_j . Note that the weight of the link between two areas does not depend only on their (normalized) correlation $r(X_i, X_j)$, but also on the “power” of the two areas, as captured by the standard deviation of the corresponding cumulative anomalies. Also, recall from the previous paragraph that this standard deviation is roughly proportional to the area's size, implying that larger areas will tend to have stronger connections. The link between two areas can be positive or negative, depending on the sign of the correlation term. Fig. 6 presents the cumulative distribution function (CDF) of the absolute correlation between the cumulative anomalies of areas for four SST

by the size of the cell in that resolution.

networks. As with the correlations of the cell-level network, there is no clear cutoff³ separating significant correlations from noise. For this reason we prefer to not prune the weaker links between areas. Instead, every pair of areas A_i and A_j is connected through a weighted link and the resulting graph is *complete*.

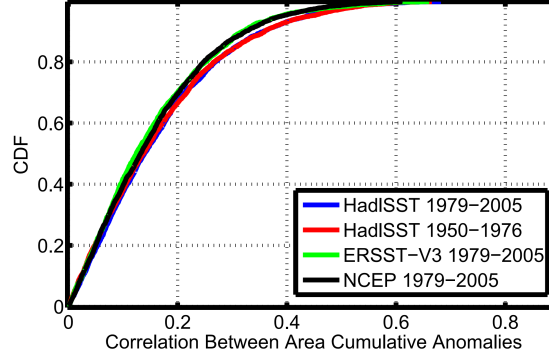


Figure 6: CDF of the absolute correlation between area cumulative anomalies for the HadISST reanalysis during the 1950-1976 and 1979-2005 periods, and for ERSST-V3 and NCEP during the 1979-2005 period

3.4 Network metrics

We now proceed to define a few network metrics that are used throughout the chapter. A climate network N is defined by a set V of areas $A_1, \dots, A_{|V|}$, representing the nodes of the network, and a set of link weights, given by Eq. 4. Because the network is a complete weighted graph, basic graph theoretic metrics that do not account for link weights (such as average degree, average path length, or clustering coefficient) are not relevant in this context.

A first representation of the network can be obtained through *link maps*. The link map of an area A_k shows the weight of the links between A_k and every other area in the network. Link maps provide a direct visualization of the correlations, positive and negative, between a given area and others in the system, often related to atmospheric teleconnection

³Imposing a threshold on the actual strength of the link (computed as the covariance between the cumulative anomalies of two areas) would be incorrect. For example, multiplying low correlations with large standard deviations can produce links of significant weight.

patterns. For instance, Fig. 7 shows link maps for the two largest areas identified in the HadISST network in the 1979-2005 period. The first area has a clear correspondence to the El Niño Southern Oscillation (ENSO); indeed, the cumulative anomaly over that area and most common indices that describe ENSO variability are highly correlated (the correlation reaches 0.94 for the Niño-3.4 index). The links of this “ENSO” area depict known teleconnections and their strength. The second largest area covers most of the tropical Indian Ocean and represents the region that is most responsive to interannual variability in the Pacific. It corresponds, broadly, to the region where significant warming is observed during peak El Niño conditions [32].

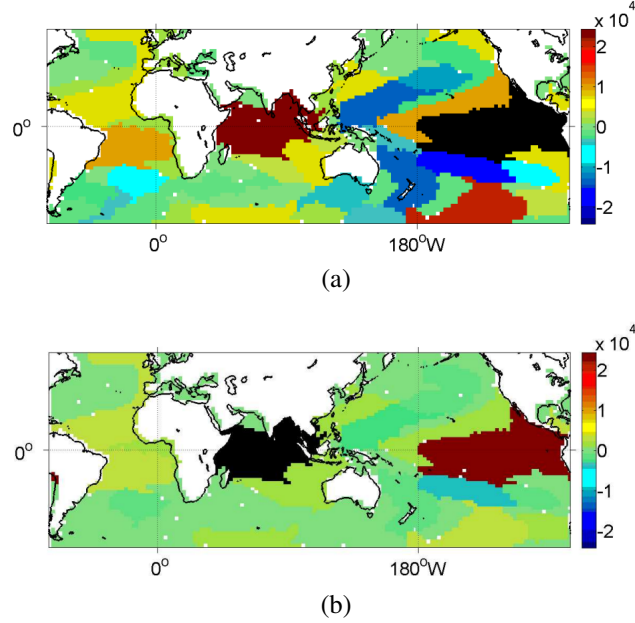


Figure 7: Link maps for two areas related to (a) ENSO and (b) the equatorial Indian Ocean in the HadISST 1979-2005 network ($\tau=0.496$). The color scale represents the weight of the link between the area shown in black and every other area in this SST network

Another metric is the *strength* of an area (also known as weighted degree), defined as the sum of the absolute link weights of that area,

$$W(A_i) = \sum_{j \neq i}^V |w(A_i, A_j)| = s(X_i) \sum_{j \neq i}^V s(X_j) |r(X_i, X_j)|. \quad (5)$$

Note that anti-correlations (negative weights) also contribute to an area’s strength. Fig. 8

shows, for example, the strength maps for two HadISST networks covering the 1950-1976 and 1979-2005 periods, respectively. Both the geographical extent of areas and their strength display differences in the two time intervals, particularly in the North Pacific sector and in the tropical Atlantic [110, 125].

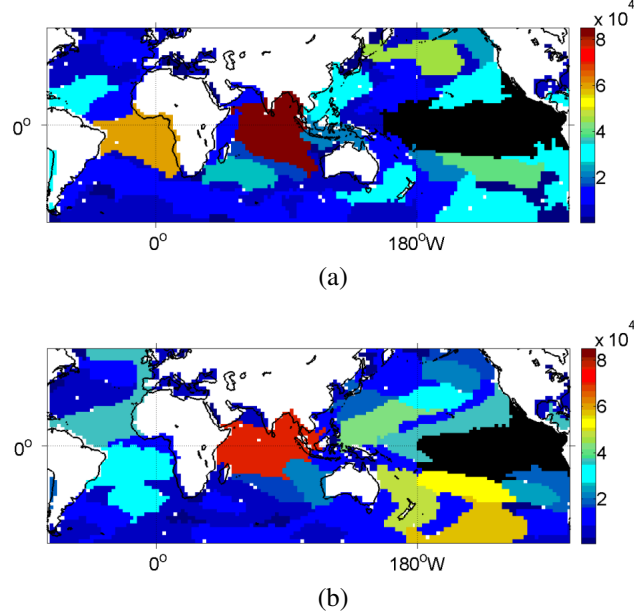


Figure 8: Strength maps for two different time periods using the HadISST data set. (a) 1950-1976 network, strength of ENSO area: 20.1×10^4 ; (b) 1979-2005 network, strength of ENSO area: 18.8×10^4

It is often useful to “peel” the nodes of a network in successive layers of increasing network significance. For weighted networks, we can do so through an iterative process referred to as *s-core* decomposition [161]. The areas of the network are first ordered in terms of their strength. In iteration-1 of the algorithm, the area with the minimum strength, say W_{min} , is removed. Then we recompute the (reduced) strength of the remaining areas, and if there is an area with lower strength than W_{min} , it is removed as well. Iteration-1 continues in this manner until there is no area with strength less than W_{min} . The areas removed in this first iteration are placed in the same layer. The algorithm then proceeds similarly with iteration-2, forming the second layer of areas. The algorithm terminates when we have removed all areas, say after K iterations. Finally, the K layers are re-labeled

as “cores” in inverse order, so that the *first order core* consists of the areas removed in the last iteration (the strongest network layer), while the *Kth order core* consists of the areas removed in the first iteration (the weakest layer). Fig. 9 shows the *top five cores* for two HadISST networks, covering 1950-1976 and 1979-2005, respectively. Again, changes in the relative role of areas are apparent in the North Pacific and in the tropical Atlantic.

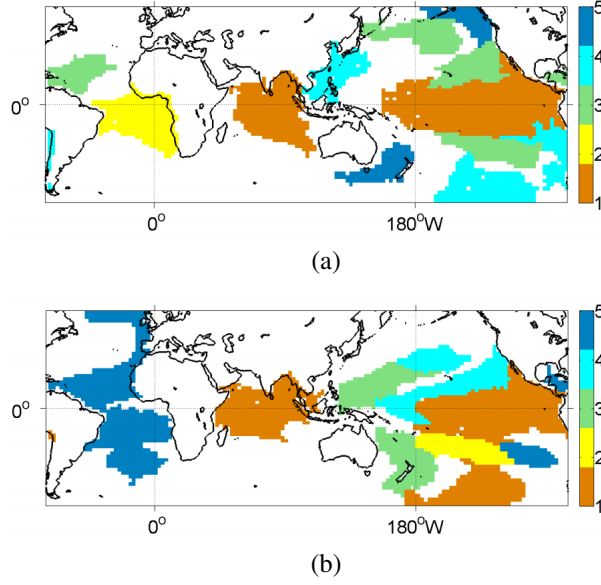


Figure 9: Color maps depicting the *top-5 order cores* for the (a) HadISST 1950-1976, and (b) HadISST 1979-2005 networks

Visual network comparisons provide insight but quantitative metrics that summarize the distance between two networks into a single number would be useful. A challenge is that the climate networks under comparison may have a different set of areas, and it is not always possible to associate an area of one network with a unique area of another network.

We rely on two quantitative metrics: the *Adjusted Rand Index* (ARI), which focuses on the similarity of two networks in terms of the identified areas, and the *Area Strength Distribution Distance*, or simply *Distance* metric, which considers the magnitude of link weights and thus area strengths.

The (non-adjusted) Rand Index is a metric that quantifies the similarity of two partitions of the same set of elements into non-overlapping subsets or “clusters” [120]. Every pair

of elements that belong to the same cluster in both partitions, or that belong to different clusters in both partitions, contributes positively to the Rand Index. Every pair of elements that belong to the same cluster in one partition but to different clusters in the other partition, contributes negatively to the Rand Index. The metric varies between 0 (complete disagreement between the two partitions) to 1 (complete agreement). A problem with the Rand Index is that two random partitions would probably give a positive value because some agreement between the two partitions may result by chance. The Adjusted Rand Index (ARI) [86, 140] ensures that the expected value of ARI in the case of random partitions is 0, while the maximum value is still 1. We refer the reader to the previous references for the ARI mathematical formula.

In the context of our method, the common set of elements is the set of grid cells, while a partition represents how cells are classified into areas (i.e., each area is a cluster of cells). Cells that do not belong to any area are assigned to an artificial cluster that we create just for computing the ARI metric. We use the ARI metric to evaluate the similarity of two networks in terms of the identified areas. This metric, however, does not consider cell anomalies and cell sizes, and so it cannot capture similarities or differences between two networks in terms of link weights, and area strengths. Two networks may have some differences in the number or spatial extent of their areas, but they can still be similar if those “ambiguously clustered” cells do not have a significant anomaly compared to their area’s anomaly. Also, two networks can have similar areas but the magnitude of their area anomalies can differ significantly, causing significant differences in link weights and thus area strengths. Further, the ARI metric cannot be used to compare data sets with different resolution because the underlying set of cells in that case would be different between the two networks.

For these reasons, together with the ARI, we rely on a distance metric that is based on the area strength distribution of the two networks. The strength of an area, in effect, summarizes the combined effect of the area’s spatial scope (which cells participate in that

area), and of the anomaly and size of those cells.

Given two networks N and N' with V and $V' \leq V$ areas, respectively, we first add $V - V'$ “virtual” areas of zero strength in network N' so that the two networks have the same number of nodes. Then, we rank the areas of each network in terms of strength, with A_i being the i 'th highest-strength area in network N . Fig. 10a shows the ranked area strength distributions for the HadISST networks covering 1950-1976 and 1979-2005 periods. The distance $d_{sd}(N, N')$ quantifies the similarity between two networks in terms of their ranked area strength distribution,

$$d_{sd}(N, N') = \sum_{i=1}^V |W(A_i) - W(A'_i)| \quad (6)$$

To normalize the previous metric, we introduce the *relative distance* $D_{sd}(N, N')$. Specifically, we construct an ensemble of randomized networks N_r with the same number of areas and link weight distribution as network N , but with random assignment of links to areas. The random variable $d_{sd}(N, N_r)$ represents the distance between N and a random network N_r , while $\overline{d_{sd}(N, N_r)}$ denotes the sample average of this distance across 100,000 such random networks. The relative distance $D_{sd}(N, N')$ is then defined as

$$D_{sd}(N, N') = \frac{d_{sd}(N, N')}{\overline{d_{sd}(N, N_r)}}. \quad (7)$$

Note that $D_{sd}(N, N')$ represents an ordered relation, from network N to N' . A relative distance close to 0 implies that N' is similar to N in terms of the allocation of link weights to areas. As the relative distance approaches 1, N' may have a similar link weight distribution with N , but the two networks differ significantly in the assignment of links to areas. The relative distance can be larger than 1 when N' 's link weight distribution is significantly different than that of N .

Two networks may be similar in terms of the identified areas (high ARI) but with large distance (high D_{sd}) if the strength of at least some areas is significantly different across the two networks (perhaps due to the magnitude of the underlying cell anomalies). In principle, it could also be that two networks have similar ranked area strength distributions

(low D_{sd}) but significant differences in the number or spatial extent of the identified areas. Consequently, the joint consideration of both metrics allows us to not only evaluate or rank pairs of networks in terms of their similarity, but also to understand which aspects of those pairs of networks are similar or different.

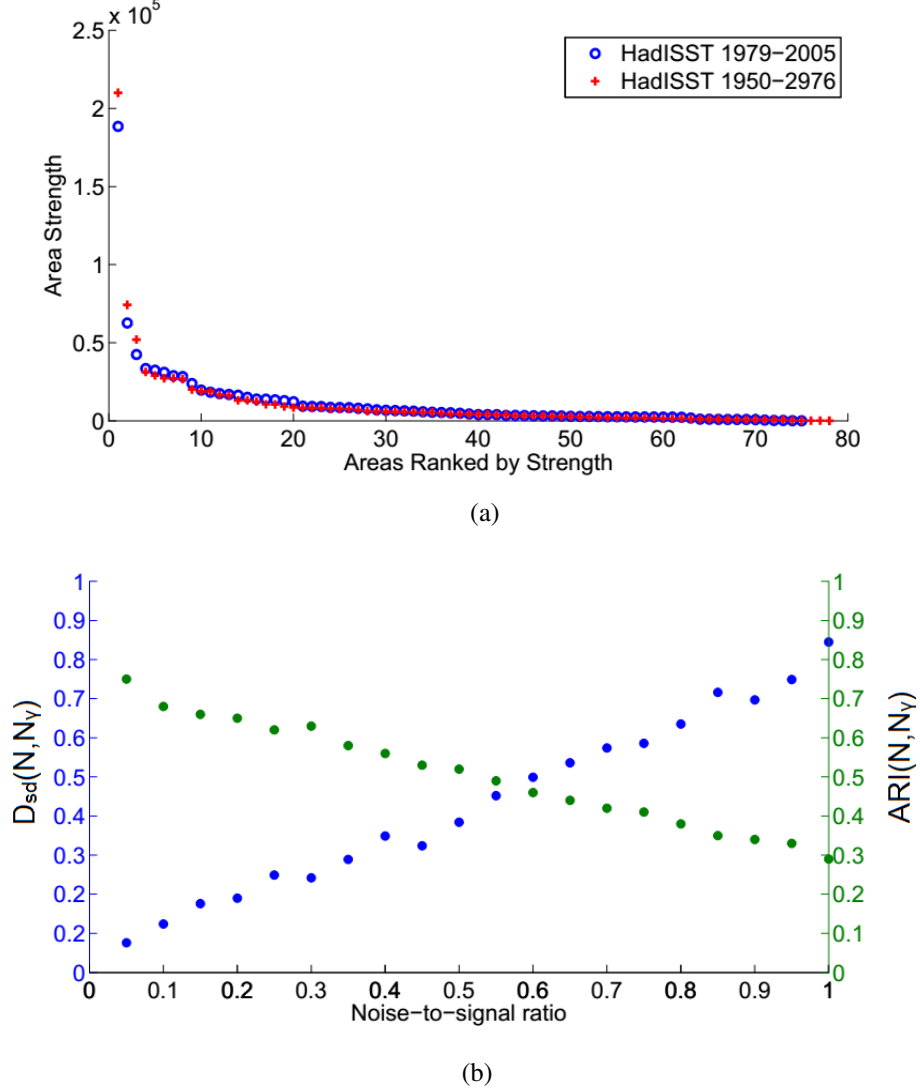


Figure 10: (a) Distribution of ranked area strengths for two networks constructed using the HadISST data set over the periods 1950-1976 and 1979-2005, respectively. (b) Distance $D_{sd}(N, N_\gamma)$ and $ARI(N, N_\gamma)$ between the HadISST 1979-2005 network and networks constructed after the addition of white Gaussian noise in the same data set

We can also map a distance $D_{sd}(N, N')$ to an amount of White Gaussian Noise (WGN) that, if added to the climate field that produced N , will result in a network with equal

distance from N . In more detail, let $s^2(x_i)$ be the sample variance of the anomaly time series $x_i(t)$ in the climate field under consideration. We construct a perturbed climate field by adding WGN with variance $\gamma s^2(x_i)$ to every $x_i(t)$, where γ is referred to as the *noise-to-signal ratio*. Then, we construct the corresponding network N_γ , and $D_{sd}(N, N_\gamma)$ is its distance from N . A given distance $D_{sd}(N, N')$ can be mapped to a noise-to-signal ratio γ when $D_{sd}(N, N') = D_{sd}(N, N_\gamma)$. Similarly, a given ARI value $\text{ARI}(N, N')$ can be mapped to noise-to-signal ratio γ such that $\text{ARI}(N, N') = \text{ARI}(N, N_\gamma)$. Fig. 10b shows how γ affects $D_{sd}(N, N_\gamma)$ and $\text{ARI}(N, N_\gamma)$ when the network N corresponds to the HadISST 1979-2005 reanalysis. As a reference point, note that a low noise magnitude, say $\gamma=0.1$, corresponds to distance $D \approx 0.12$ and $\text{ARI} \approx 0.68$.

Finally, we emphasize that the ARI and D_{sd} metrics focus on the global scale. Even if two networks are quite similar according to these two metrics, meaningful differences at the local scale of individual areas may still exist. The study of regional climate effects may require an adaptation of these metrics.

3.5 Robustness analysis

Analyzing climate data poses many challenges: measurements provide only partial geographical and temporal coverage, while the collected data are subject to instrumental biases and errors both random and systematic. Greater uncertainties exist in general circulation model outputs: climate simulations are dependent on modeling assumptions, complex parameterizations and implementation errors. An important question for any method that identifies topological properties of climate fields is whether it is robust to small perturbations in the input data, the method parameters, or in the assumptions the method is based on. If so, the method can provide useful information on the climate system despite uncertainties of various types. In this section, we examine the sensitivity of the inferred networks to deviations in the input data, the parameter τ , and certain methodological choices. In all cases we quantify sensitivity by computing the D_{sd} and ARI metrics from the original

network to each of the perturbed networks.

3.5.1 Robustness to additive white Gaussian noise

As described in Section 3.4, a simple way to perturb the input data is to add white Gaussian noise to the original climate field time series. The magnitude of the noise is controlled by the *noise-to-signal ratio* γ . The distance D_{sd} and ARI from the original network N to the “noisy” networks N_γ are shown in Fig. 10b for the HadISST reanalysis over 1979-2005. To visually illustrate how noise affects the identified areas, and in particular their strength, Fig. 11 presents strength maps for two values of γ ; the area strengths should be compared with Fig. 8b. Although some differences exist, the ENSO area strength is comparable to that of the original network, and the hierarchy (in terms of strength) in the three basins is conserved.

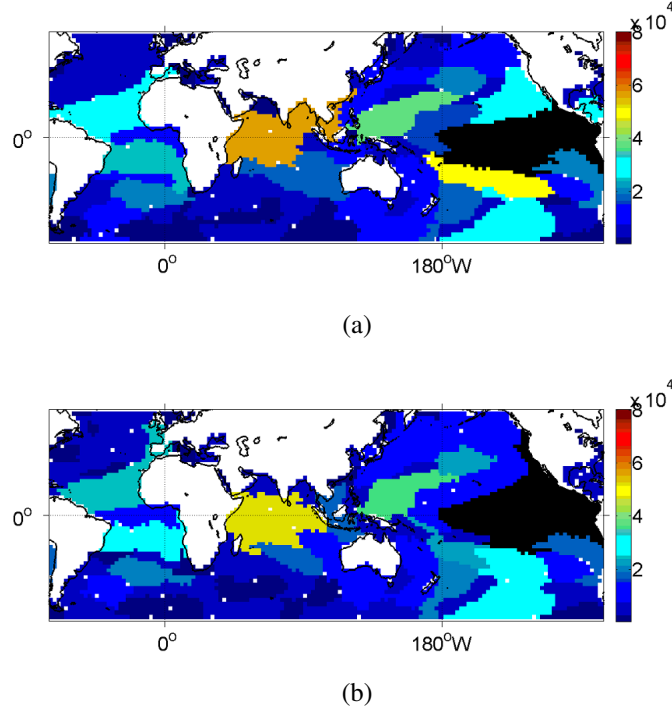


Figure 11: Strength maps for two perturbations of the HadISST 1979-2005 data set using white Gaussian noise. (a) $\gamma=0.05$, strength of ENSO area: 18.0×10^4 . (b) $\gamma=0.10$, strength of ENSO area: 19.1×10^4

3.5.2 Robustness to the resolution of the input data set

All data sets compared in this chapter have been spatially interpolated to the lowest common resolution. Here we investigate the robustness of the identified network to the resolution of the input data set. To do so, consider the HadISST reanalysis over the 1979-2005 period and compare the network discussed so far, constructed using data interpolated on a $2^\circ\text{lat} \times 2.5^\circ\text{lon}$ grid, with two networks based on a lower ($4^\circ\text{lat} \times 4^\circ\text{lon}$) and a higher ($1^\circ\text{lat} \times 2^\circ\text{lon}$) resolution realization of the same reanalysis. Fig. 12 shows strength maps for the two new networks. As we lower the resolution the total number of areas decreases, and the areas immediately surrounding the ENSO-related area get weaker. Nonetheless, the hierarchy of area strengths in the three basins is preserved, and differences are small, as quantified by the distance metric. The distance from the default to the high resolution network is $D_{sd}(N, N')=0.10$ ($\gamma=0.07$). The distance from the default to the low resolution network is $D_{sd}(N, N')=0.11$ ($\gamma=0.10$). As previously mentioned, the ARI cannot be used to compare data sets with different spatial resolution.

3.5.3 Robustness to the selection of τ

Recall that the parameter τ represents the threshold for the minimum average pair-wise correlation between cells of the same area. Even though we provide a heuristic (see Section 3.8) for the selection of τ , which depends on the given data set, it is important to know whether small deviations in τ have a major effect on the constructed networks.

Considering again the HadISST 1979-2005 reanalysis, Fig. 13 presents the relative distance and ARI from the original network N constructed using $\tau=0.496$ (it corresponds to a significance level $\alpha = .1\%$), to networks N_τ constructed using different τ values. We vary τ by $\pm 10\%$, in the range 0.45–0.55. This corresponds to a large change, roughly an order of magnitude, in the underlying significance level α .

Fig. 14 visualizes strength maps for the two extreme values of τ in the previous range. While some noticeable differences exist, the overall area structure appears robust to the

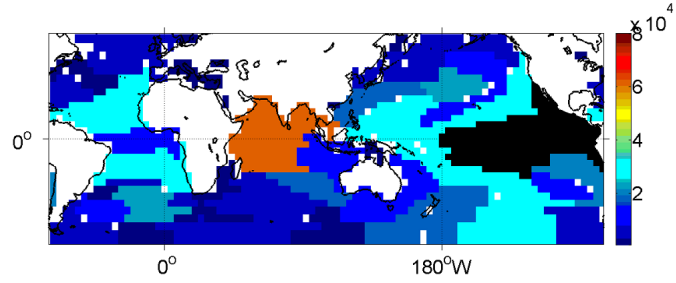
choice of τ . By increasing τ , we increase the required degree of homogeneity within an area, and therefore the resulting network will be more fragmented, with more areas of smaller size and lower strength, and vice versa for decreasing τ .

3.5.4 Robustness to the selection of the correlation metric

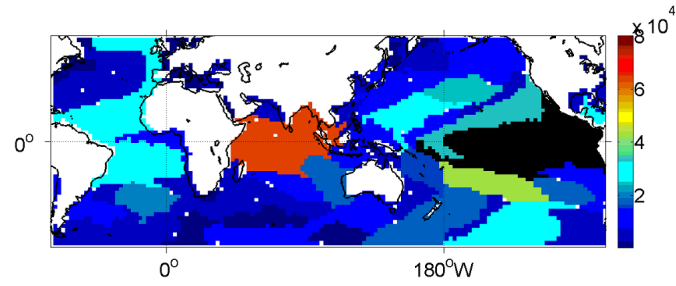
The input to the network construction process is a matrix of correlation values between all pairs of cells. So far, we have relied on Pearson's correlation coefficient, which is a linear dependence measure between two random variables. Any other correlation metric could be used instead. To verify that the properties of the resulting network do not depend strongly on the selected correlation metric, we use here the non-parametric *Spearman's rank coefficient* to compute cell-level correlations.

Fig. 15 shows the strength map for the HadISST 1979-2005 network using Spearman's correlation metric. Again, while small changes are apparent, the size and shape of the major areas and their relative strength are unaltered. $D_{sd}(N, N')=0.08$ and $\text{ARI}(N, N')=0.76$, where N is the network shown in Fig. 8b; both metrics correspond to $\gamma=0.05$.

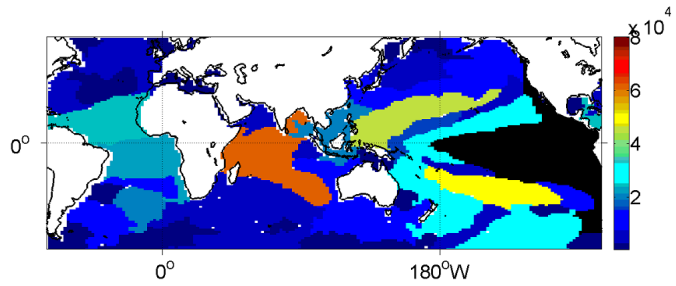
We have performed similar robustness tests using precipitation data obtaining comparable results.



(a)

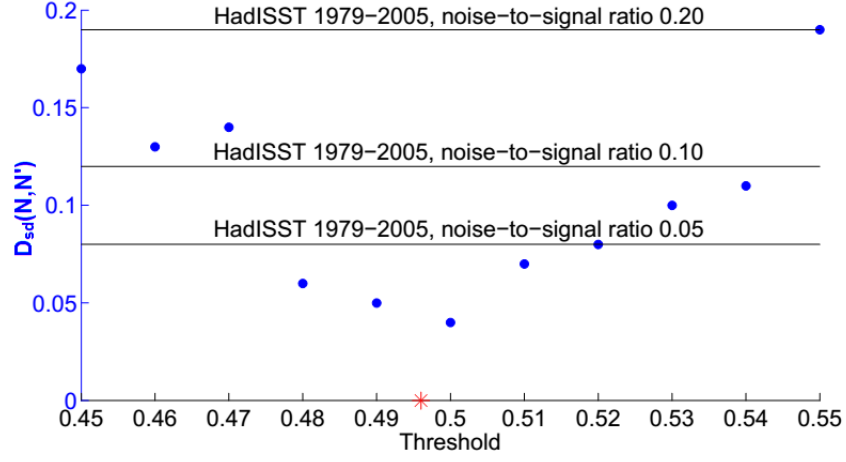


(b)

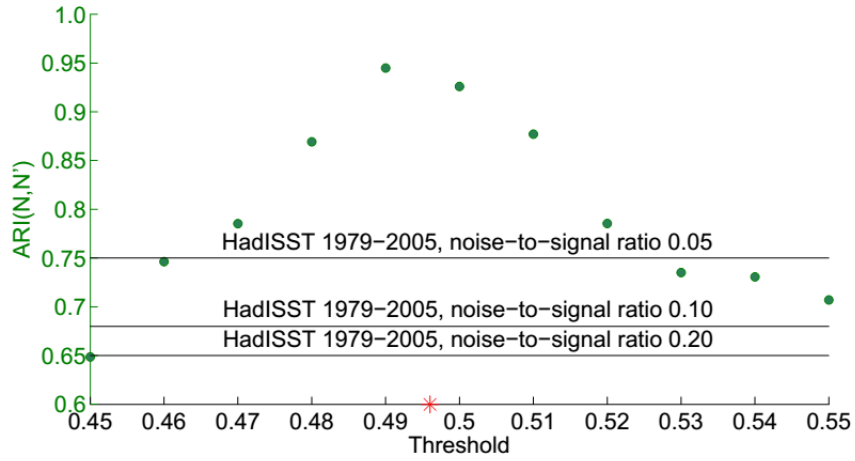


(c)

Figure 12: Strength maps for the HadISST 1979-2005 network at three different resolutions. (a) Low resolution network, ($4^{\circ}lat \times 4^{\circ}lon$), strength of ENSO area: 18.2×10^4 . (b) Default resolution network, ($2^{\circ}lat \times 2.5^{\circ}lon$), strength of ENSO area: 18.8×10^4 . (c) High resolution network, ($1^{\circ}lat \times 2^{\circ}lon$), strength of ENSO area: 18.2×10^4

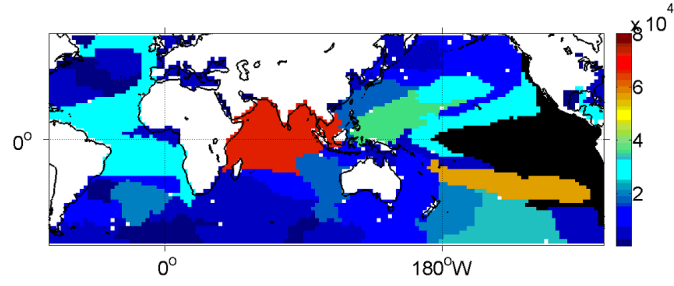


(a)

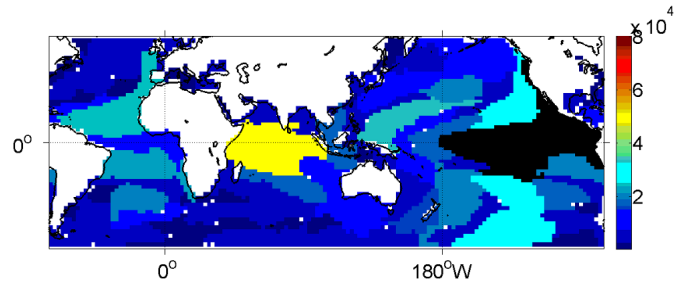


(b)

Figure 13: (a) Distance D_{sd} and (b) ARI from the original HadISST 1979-2005 network (marked with an asterisk in the x-axis, $\tau=0.496$) to networks constructed with different values of τ . The black horizontal lines correspond to the distance $D_{sd}(N, N_\gamma)$ and $ARI(N, N_\gamma)$



(a)



(b)

Figure 14: Strength maps for the HadISST 1979-2005 network using two values of the parameter τ . The “default” value is $\tau=0.496$, corresponding to $\alpha=.1\%$ (see Section 3.8).
 (a) $\tau=0.45$, strength of ENSO area: 18.7×10^4 . (b) $\tau=0.55$, strength of ENSO area: 18.6×10^4

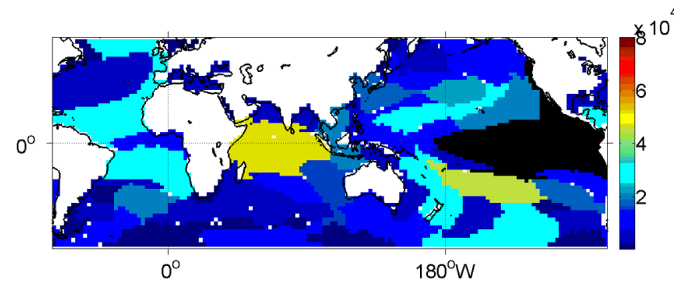


Figure 15: Strength map for the HadISST 1979-2005 network using Spearman's correlation; strength of ENSO area: 18.5×10^4

3.6 Applications

We now apply the proposed method to the climate data sets described in Section 3.2 to illustrate that network analysis can be successfully used to compare data sets and to validate model representations of major climate areas and their connections. We proceed by constructing networks for three different SST reanalyses and two precipitation data sets. We then examine the relation between two different climate fields (SST and precipitation) introducing a *regression of networks* technique. Finally, we analyze the network structure of the SST fields from two models participating in CMIP5.

3.6.1 Comparison of SST networks

Here we investigate the network properties and metrics for three SST reanalyses focusing on the 1979-2005 period. Two of them, HadISST and ERSST-V3, use statistical methods to fill sparse SST observations; HadISST implements a reduced space optimal interpolation (RSOI) technique, while ERSST-V3 adopts a method based on empirical orthogonal function (EOF) projections. NCEP/NCAR uses the Global Sea Ice and Sea Surface Temperatures (GISST2.2) from the U.K. Meteorological Office until late 1981 and the NCEP Optimal Interpolation (OI) SST analysis from November 1981 onward. The GISST2.2 is based on empirical orthogonal function (EOF) reconstructions [87]. The OI SST analysis technique combines in situ and satellite-derived SST data [123]. To minimize the possibility of artificial trends, and the bias introduced by merging different data sets, GISST data are modified to include an EOF expansion based on the IO analysis from January 1982 to December 1993.

In Fig. 16, we quantify the differences between the three reanalyses showing correlation maps between the detrended DJF SST anomaly time series for HadISST and ERSST-V3, HadISST and NCEP, and ERSST-V3 and NCEP. The patterns that emerge in the all correlation maps are similar. Correlations are generally higher than 0.9 in the equatorial Pacific, due to the almost cloud free sky and to the in-situ coverage provided since the mid 80s'

first by the Tropical Ocean Global Atmosphere (TOGA) program, and then by the Tropical Atmosphere Ocean (TAO)/Triangle Trans-Ocean Buoy Network (TAO/TRITON) program [166]. Good agreement between reanalyses is also found in the north-east Pacific, in the tropical Atlantic and in the Indian and Pacific Oceans between 10° S and 30° S. Correlations decrease to approximately 0.7 in the equatorial Indian Ocean and around Indonesia, where cloud coverage limits satellite retrievals, and reach values as small as 0.2-0.3 in the Labrador Sea, close to the Bering Strait and south of 40° S, particularly in the Atlantic and Indian sectors, due to persistent clouds and poor availability of in-situ data. North of 60° N and south of 60° S the presence of inadequately sampled sea-ice and intense cloud coverage reduce even further the correlations, that attain non-significant values almost everywhere. At those latitudes any comparison between those reanalyses and their resulting networks is meaningless given that it would not possible to identify a reference data set.

The strength maps constructed using these data sets show differences in all basins, and suggest that the network analysis performed allows for capturing more subtle properties than correlation maps (Fig. 17). To begin with the strongest area, corresponding to ENSO, we notice that it has a similar shape in HadISST and NCEP, but it extends further to the west in ERSST-V3. Its strength is about 10% higher in NCEP compared to the other two reanalyses. In HadISST, the equatorial Indian Ocean appears as the second strongest area, followed by areas surrounding the ENSO region in the tropical Pacific and by the tropical Atlantic. In ERSST-V3 the area comprising the equatorial Indian Ocean has shape and size analogous to HadISST, but 30% weaker, and it is closer in strength to the area covering the warm-pool in the western tropical Pacific. Also the areas comprising the tropical Atlantic are slightly weaker than in the other two data sets. HadISST and ERSST-V3 display a similar strength hierarchy, with the Pacific Ocean being the basin with the strongest (ENSO-like) area, followed by the Indian, and finally by the Atlantic Ocean. In NCEP all tropical areas (except the area corresponding to the ENSO region) have similar strength and the hierarchy between Indian and Atlantic Oceans is inverted. Also, the equatorial Indian

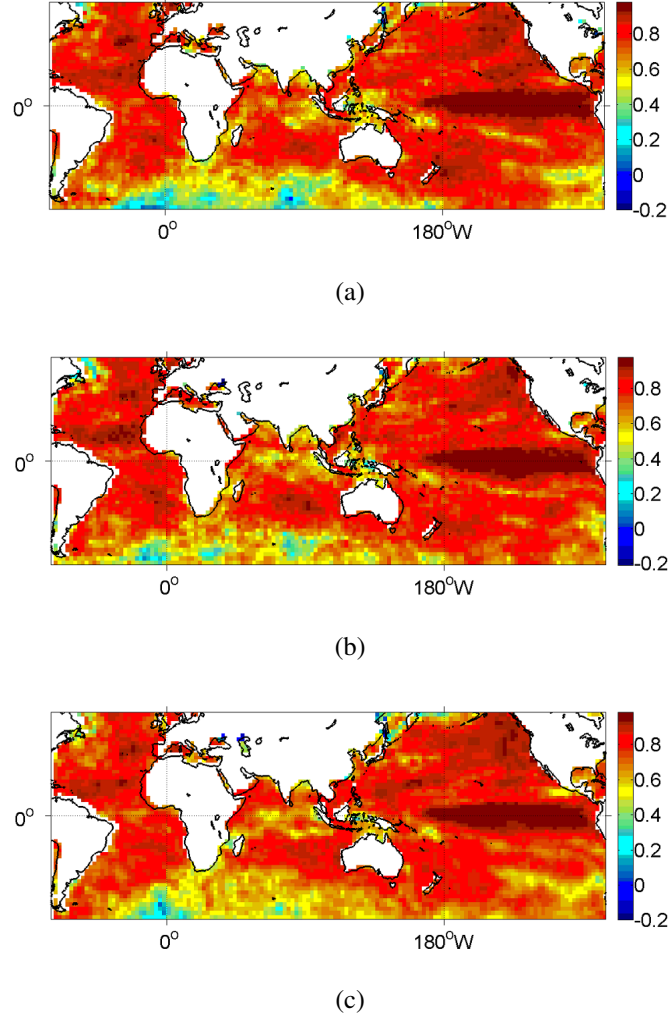


Figure 16: Pearson correlation maps between the SST anomaly time series in all pairs of three reanalyses data sets over the 1979-2005 period in boreal winter (DJF). Correlations between (a) HadISST and ERSST-V3; (b) HadISST and NCEP; (c) NCEP and ERSST-V3

Ocean appears subdivided in several small areas.

Differences in strength maps are also reflected in the *s-core* decomposition (Fig. 18) and in the links between the ENSO-related areas and other areas in the network (Fig. 19). In HadISST and ERSST-V3, the *first order core* is located in the tropical and equatorial Pacific and Indian Ocean, while in NCEP it is limited to the Pacific. As a consequence the strength of the link between the ENSO-related area and the Indian Ocean is much stronger in the first two reanalyses than in NCEP. In HadISST, the ENSO-related and Indian Ocean areas are

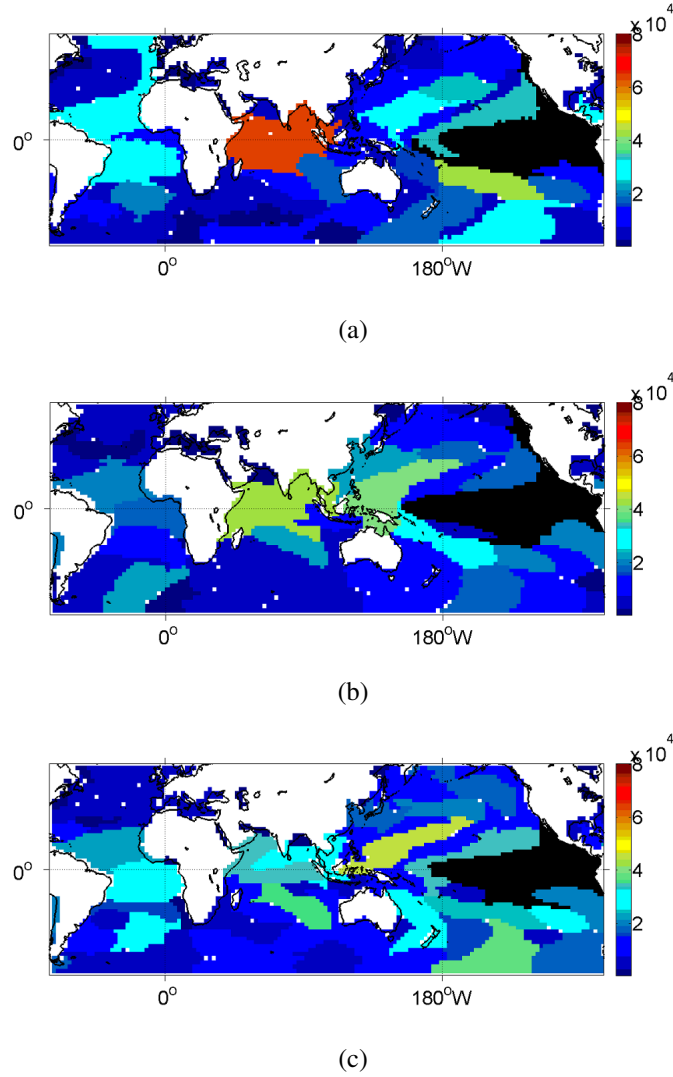


Figure 17: Strength maps for networks constructed based on (a) HadISST (ENSO area strength 18.8×10^4); (b) ERSST-V3 (ENSO area strength 17.6×10^4); (c) NCEP (ENSO area strength 21.0×10^4). In all networks the period considered is 1979-2005

separated by regions of higher order in the western Pacific, organized in the characteristic “horse-shoe” pattern. In the other two reanalyses the *first order core* extends along the whole Pacific equatorial band and includes the horse-shoe areas. In correspondence, the links between the ENSO-like and the western Pacific areas are, in absolute value, weaker than the link between ENSO and the Indian Ocean in HadISST, but comparable in ERSST-V3. NCEP shows significantly weaker links overall, but the highest link weights are found between ENSO and the western Pacific.

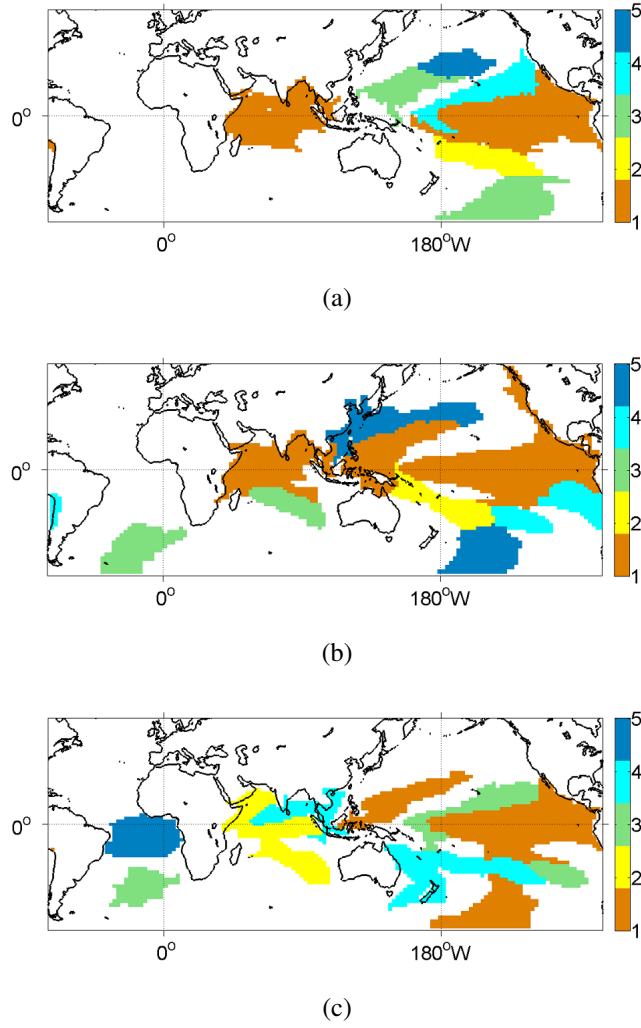


Figure 18: *Top-5 order cores* in (a) HadISST; (b) ERSST-V3; (c) NCEP. The period considered is 1979-2005 in all cases

To conclude the comparison of different SST reanalyses, we measure the distance and ARI values from HadISST to the other two networks. The distance from HadISST to ERSST-V3 is small, $D_{sd}(N, N')=0.16$, mapped to a *noise-to-signal* ratio $\gamma=0.15$. The strongest areas show indeed a good correspondence in strength and size in the two data sets, even if the shape of the ENSO-related areas differ. The distance from HadISST to NCEP, $D_{sd}(N, N')=0.29$ with $\gamma=0.35$, is greater, as expected from the previous figures, given that all areas except of the ENSO-related one appear significantly weaker, while the ENSO area is stronger than in HadISST. NCEP is also *penalized* because of the differences, compared

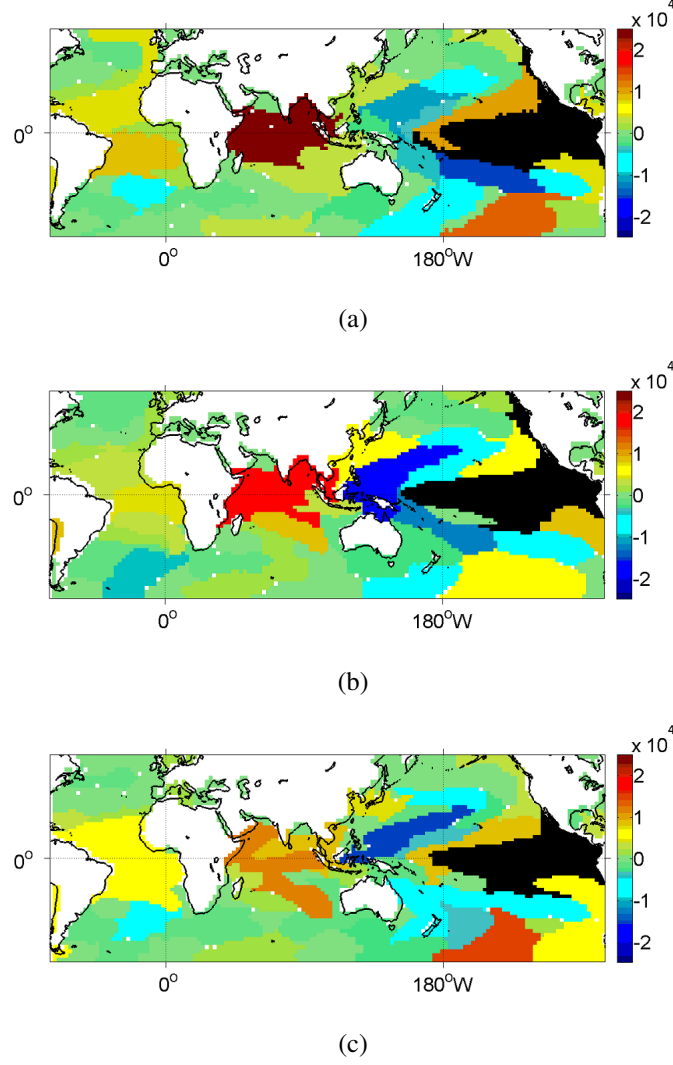


Figure 19: Links between the ENSO-like area shown in black and all other areas in the three reanalyses. (a) HadISST, (b) ERSST-V3 and (c) NCEP networks

to HadISST, in the strength (and size) of areas over the Indian Ocean and in the horse-shoe pattern. Recall that D_{sd} compares areas based on their strength ranking, independent on their geographical location. In this respect, the two strongest areas represented by ENSO and Indian Ocean in HadISST are replaced by ENSO and the North Pacific extension of the horse-shoe region in NCEP. The ARI metric, on the other hand, ranks NCEP closer to HadISST than ERSST-V3 (ARI=0.59 for NCEP and ARI=0.54 for ERSST-V3, mapped to $\gamma=0.35$ and 0.45, respectively). The shape of the ENSO-related area and of areas in the tropical Atlantic and south of 30° S are indeed in better agreement between HadISST and

NCEP, despite having different strengths.

The previous discussion illustrates that D_{sd} and ARI should be considered jointly, as they provide complementary information about the similarity and differences between two networks.

3.6.2 Network changes over time

Network analysis can also be a powerful tool to detect and quantify climate shifts. The insights that network analysis can offer, compared to more traditional time series analysis methods, are related to the detection of changes in network metrics that are associated with specific climate modes of variability, regional or global. Topological changes may include addition or removal of areas, significant fluctuations in the weight of existing links (strengthening and weakening of teleconnections), or variations in the relative significance of different areas, quantified by the area strength distribution. For instance, Tsonis and co-authors have built a network of four interacting nodes using the major climate indices, the North Atlantic Oscillation (NAO), ENSO, the North Pacific Oscillation (NPO) and the Pacific Decadal Oscillation (PDO), and suggested that those climate modes of variability tend to synchronize with a certain coupling strength [154]. Climate shifts, including the one recorded in the north Pacific around 1977 [110], could result from changes in such coupling strength.

Here we compare the climate networks constructed on the HadISST data set over the periods 1950-1976 and 1979-2005 to illustrate that the proposed methodology may also provide insights into the detection of climate shifts. Instead of simply comparing different periods, it is possible to use a sliding window in the network inference process to detect significant changes or shifts without prior knowledge; we will explore this possibility in future work.

Strength maps for the two networks were shown in Fig. 8, while the *top-5 order cores*

were shown in Fig. 9. The links from the ENSO-related area and from the equatorial Indian Ocean during the 1950-1976 period are presented in Fig. 20, and they can be compared with Fig. 7. When the 1979-2005 period is compared to the earlier period, we note a substantial strength decrease for the area covering the south tropical Atlantic and a significant weaker link between this area and ENSO. This suggests an alteration in the Pacific-Atlantic connection, which indeed has been recently pointed out by [125] and may be linked to the Atlantic warming [102]. Additionally, there is a change in the sign of the link weight between the ENSO area and the area off the coast of Alaska in the north Pacific, which is related to the change in sign of the PDO in 1976-1977 [110, 78].

Despite those differences, the distance from the 1979-2005 HadISST network to the 1950-1976 network is less than the distance from the former to any of the other reanalyses investigated earlier: $D_{sd}(N, N')=0.13$ with noise $\gamma=0.10$. The ARI, on the other hand, is 0.55 ($\gamma=0.40$). The ARI value reflects, predominantly, the changes in shape and size of the ENSO-related areas and of the areas over the North Atlantic and North Pacific.

3.6.3 Comparison of precipitation networks

One of the advantages of the proposed methodology is its applicability, without modifications, to any climate variable. As an example, in the following we focus on precipitation, chosen for having statistical characteristics very different from SST due to its intermittency. We investigate the network structure of the CPC Merged Analysis of Precipitation (CMAP) [172] and ERA-Interim reanalysis [46]. Both data sets are available from 1979 onward. CMAP provides gridded, monthly averaged precipitation rates obtained from satellite estimates. ERA-Interim is the outcome of a state-of-the-art data assimilative model that assimilates a broad set of observations, including satellite data, every 12 hours. As in the case of SSTs, we present the precipitation networks focusing on boreal winter (December to January) based on detrended anomalies from 1979 to 2005. Fig. 21 shows the map of area strengths for both data sets, Fig. 22 presents the *top-5 order cores*, while Fig. 23 depicts

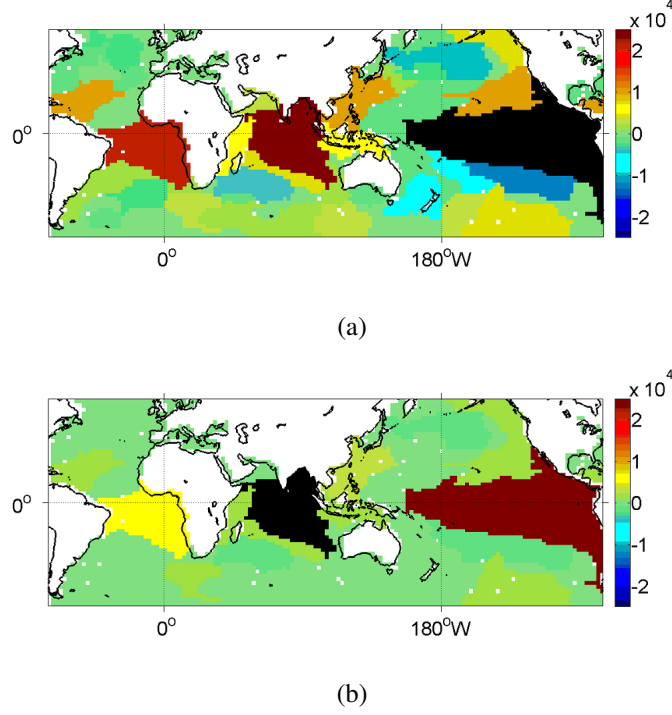
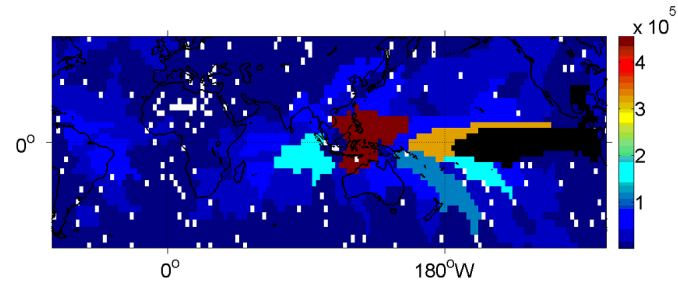


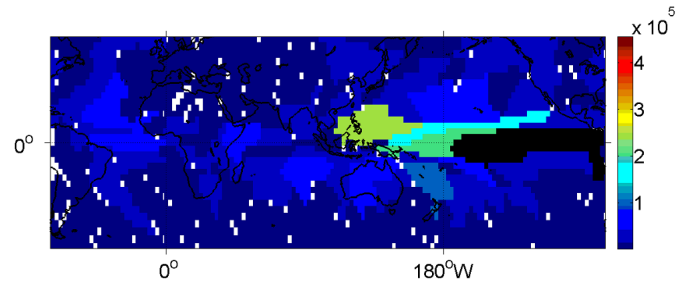
Figure 20: Links for the HadISST network over 1950 - 1976 from the (a) ENSO-related area, and (b) the equatorial Indian Ocean area (in black in the two panels)

links from the strongest area in the two networks.

The precipitation network is, not surprisingly, characterized by smaller areas, compared to SSTs. Precipitation time series are indeed highly intermittent, resulting in weaker correlations between grid cells. The areas with the highest strength are concentrated in the tropics, where deep convection takes place. The strongest area is located in the equatorial Pacific in correspondence with the center of action of ENSO. In CMAP, this area is linked with strong negative correlation to the area covering the warm-pool region, and together they represent the *first order core* of this network. The *second order core* covers the eastern part of the Indian Ocean and eastern portion of the South Pacific Convergence Zone (SPCZ). Both those regions are strongly affected by the shift in convection associated with ENSO events. In the reanalysis, the warm-pool area extends predominantly into the northern hemisphere, and its strength and size, as well as the weight of its link with the ENSO-related area, are reduced. Additionally, the Indian Ocean is subdivided in small

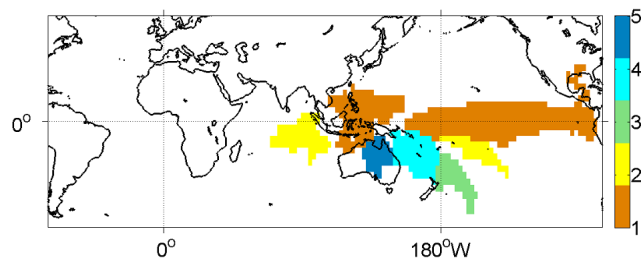


(a)

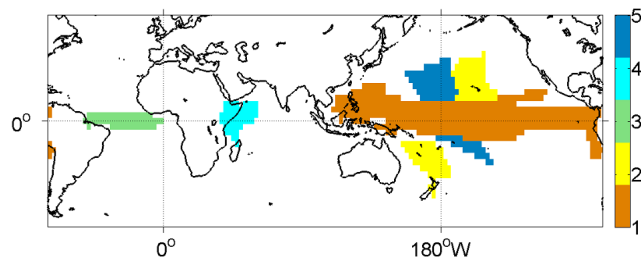


(b)

Figure 21: Precipitation networks. Area strength map in (a) CMAP (equatorial Pacific area strength 49.4×10^4), and (b) ERA-Interim (equatorial area strength 41.0×10^4)



(a)



(b)

Figure 22: *Top-5 order cores* in (a) CMAP, and (b) ERA-Interim

areas all of negligible strength, similarly to what seen for NCEP SSTs, indicating that the atmospheric teleconnection between ENSO and the eastern Indian Ocean that causes a shift in convective activity over the Indian basin (see e.g. [98, 27]) is not correctly captured by ERA-Interim. The *s-core* decomposition does not include in the *second order core* any area in the Indian Ocean, but is limited to two areas to the north and to the south of the ENSO-related one.

The distance from the CMAP network to the ERA-Interim network is $D_{sd}(N, N')=0.21$, with $\gamma=0.25$, while the ARI value is 0.49, with $\gamma=0.45$. These values reflect larger differences compared to the SST networks we presented earlier, but precipitation is known to be one of the most difficult fields to model, even when assimilating all available data, due to biases associated with the cloud formation and convective parameterization schemes [3]. In particular D_{sd} is affected by the significant difference in the strength and size of the area over the warm-pool, and of the one between the ENSO-related area and the warm-pool, while the ARI is affected by the difference in the partitions over the warm-pool and most of the Atlantic basin.

3.6.4 Regression between networks

So far we have shown applications of network analysis considering one climate variable at a time. In climate science it is often useful to visualize the relations between two or more variables to understand, for example, how changes in sea surface temperatures may impact rainfall. A simple statistical tool that highlights such relations is provided by regression analysis. Here we apply a similar approach using climate networks.

Consider two climate networks N_x and N_y , constructed using variables $\mathbf{x}(\mathbf{t})$ and $\mathbf{y}(\mathbf{t})$, respectively. The relation between an area of N_x and the areas of N_y can be quantified based on the cumulative anomaly of each area, using the earlier link weight definition (see Eq. 4). Similarly, a link map for an area $A_i \in V_x$ can be constructed based on the link weights between the area A_i and all areas $A_j \in V_y$.

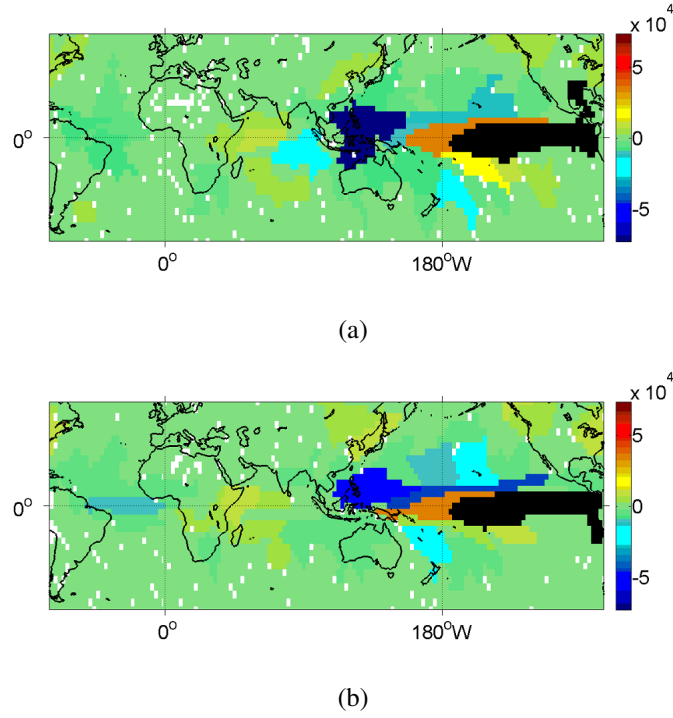


Figure 23: Link maps from the strongest area (in black) for the two precipitation reanalysis data sets. (a) CMAP; (b) ERA Interim

For instance, we construct a network linking the area that corresponds to ENSO in the HadISST reanalysis to the areas of the CMAP precipitation network for the period 1979-2005 in boreal winter. Both networks are dominated by the ENSO area and it is expected that this exercise will portrait the ENSO teleconnection patterns. Results are shown in Fig. 24. The *regression* of the rainfall network onto the ENSO-related area in the SST reanalysis visualizes the well known shift of convective activity from the warm-pool into the central and eastern equatorial Pacific during El Niño. For positive ENSO episodes, negative precipitation anomalies concentrate in the warm-pool and extend to the SPCZ and the eastern Indian Ocean. Weak, positive correlations between SST anomalies in the equatorial Pacific and precipitation are seen over the western Indian Ocean and east Africa, part of China, the Gulf of Alaska and the north-east USA. This approach is only moderately useful on reanalysis or observational data, where known indices can be used to perform regressions without the need of constructing a network. Its extension to model outputs,

however, is advantageous compared to traditional methods, because it does not require any ad-hoc index definition, but relies on areas objectively identified by the proposed network algorithm.

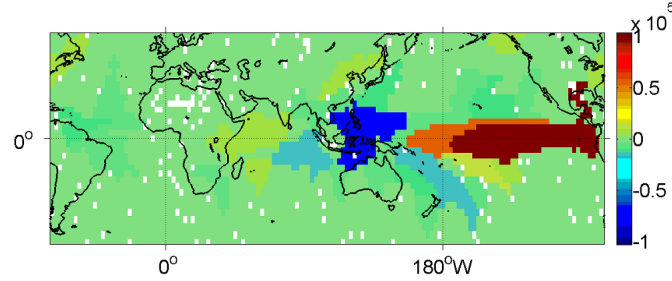


Figure 24: Link maps from the ENSO-like area in HadISST data set to all areas in the CMAP data set, considering the 1979-2005 period. Values greater than $|1 \times 10^4|$ are saturated

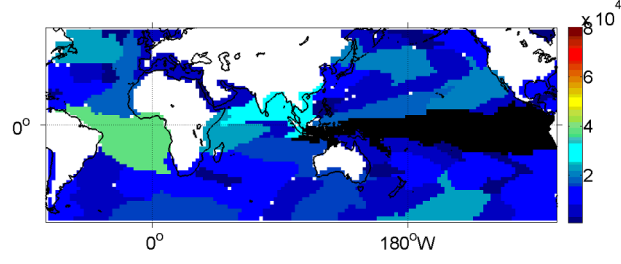
3.6.5 CMIP5 SST networks

We now compare the HadISST network with networks constructed using SST anomalies from two coupled models participating in CMIP5. Our goal is to exemplify the information that our methodology can provide when applied to model outputs. We do not aim at providing an exhaustive evaluation of the model performances, which would be beyond the scope of this chapter. We analyze the SST fields of two members of the CMIP5 historical ensemble from the GISS-E2H and HadCM3 models over the period 1979-2005. Historical runs aim at reproducing the observed climate from 1850 to 2005 including all forcings. We show strength maps (Fig. 25), *top-5 order cores* (Fig. 26), and link maps for the area that is related to ENSO (Fig. 27).

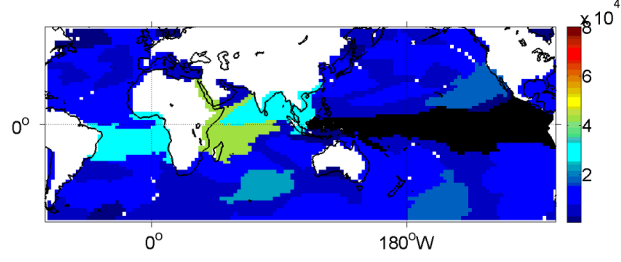
In all model integrations the ENSO-like area extends too far west into the warm-pool region, and is too narrow in the simulated width, in agreement with the recent analysis by [178]. The warm-pool is therefore not represented as an independent area anticorrelated to the ENSO-like one. In the GISS-E2H model the strength of the ENSO area is underestimated compared to the reanalyses (see Fig. 17a), but the overall size of the area is larger

than observed. Both the extent and strength of the Indian Ocean area around the equator and of the areas forming the horse-shoe pattern are reduced with respect to HadISST. Links in GISS-E2H are overall weaker than in the reanalysis (see Fig. 19a), the role of the Atlantic is slightly overestimated, and the high negative correlations between the ENSO region and the areas forming the horse-shoe patterns are not captured. In HadCM3, on the other hand, the strength of the ENSO area is comparable or greater than in the observations. In this model, areas are more numerous and fragmented than in the reanalysis, and in several cases confined within narrow latitudinal bands. This bias may result from too weak meridional currents and/or weak trade wind across all latitudes, as suggested by [179]. HadCM3 shows also erroneously strong links between the modeled ENSO area and the Southern Ocean, particularly in the Pacific and Indian sectors, as evident in the *s-core* decomposition and link maps. The link strengths in HadCM3 are closer to the observed, but some areas in the southern hemisphere play a key role, unrealistically.

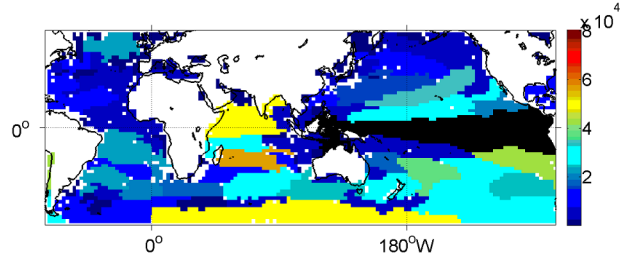
To conclude this comparison we present the distance from the HadISST reanalysis to those two models, and the corresponding ARI values. Table 2 summarizes this comparison. $D_{sd}(N, N')$ from HadISST to the two GISS-E2H integrations is 0.29 and 0.37, with $\gamma=0.35$ and $\gamma=0.45$, respectively. $D_{sd}(N, N')$ from HadISST to the two HadCM3 runs is 0.56 and 0.35, with $\gamma=0.70$ and $\gamma=0.40$. One of the GISS member networks displays a significantly smaller distance from HadISST than both networks build on the HadCM3 runs. This is due to the fact that in all networks considered the ENSO-like area overpowers all others in terms of strength and, furthermore, there exist a few other strong areas (areas that are weaker than the ENSO-related one by less than one order of magnitude). Focusing on the extent of the areas in the GISS member with smaller D_{sd} we observe striking differences relative to the base HadISST network: the GISS model is unable to reproduce the horse-shoe pattern, and it splits the tropical Indian Ocean in two areas. However, it reproduces quite well the overall size of most areas, and the strength of the largest two in the tropics, despite inverting the relative strengths of the Indian Ocean and of the south tropical Atlantic. The south tropical



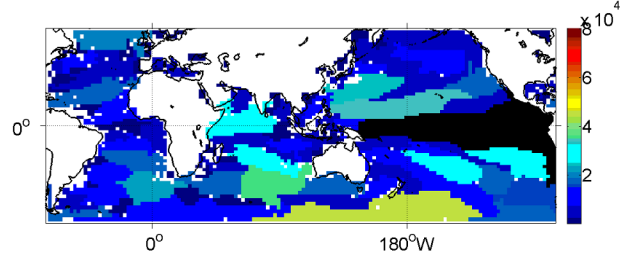
(a)



(b)



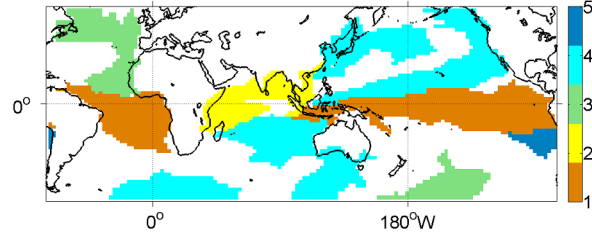
(c)



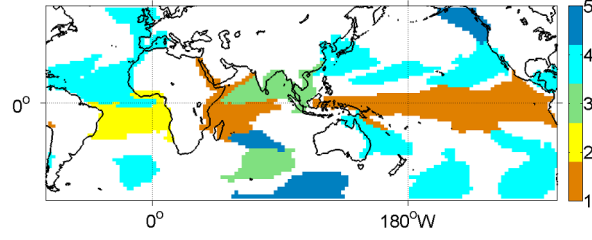
(d)

Figure 25: Strength maps for two members of the GISS-E2H and HadCM3 “historical” ensemble. (a) GISS-E2H run 1 (ENSO area strength 9.8×10^4); (b) GISS-E2H run 2 (ENSO area strength 10.0×10^4); (c) HadCM3 run 1 (ENSO area strength 23.3×10^4) and (d) HadCM3 run 2 (ENSO area strength 16.9×10^4)

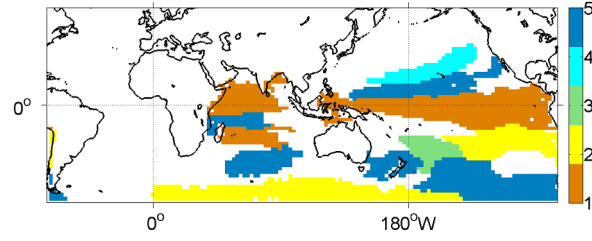
Atlantic area in GISS and the Indian Ocean one in HadISST have comparable size and strength, and D_{sd} cannot account for their different location. The HadCM3 networks, on



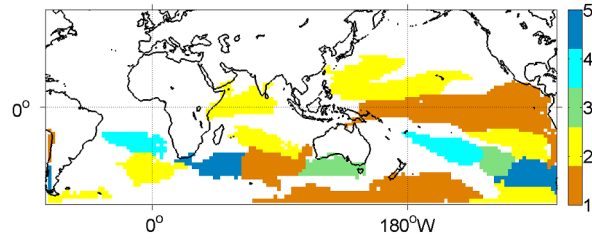
(a)



(b)



(c)



(d)

Figure 26: *Top-5 order cores* identified in the SST anomaly networks for (a-b) two GISS-E2H ensemble members and (c-d) two HadCM3 integrations

the other hand, are too fragmented and are characterized by unrealistically strong areas in the Southern Ocean, and are penalized by D_{sd} for not capturing properly the size of the strongest areas. The ARI values are 0.46 and 0.48 for the two GISS members, and 0.43 and

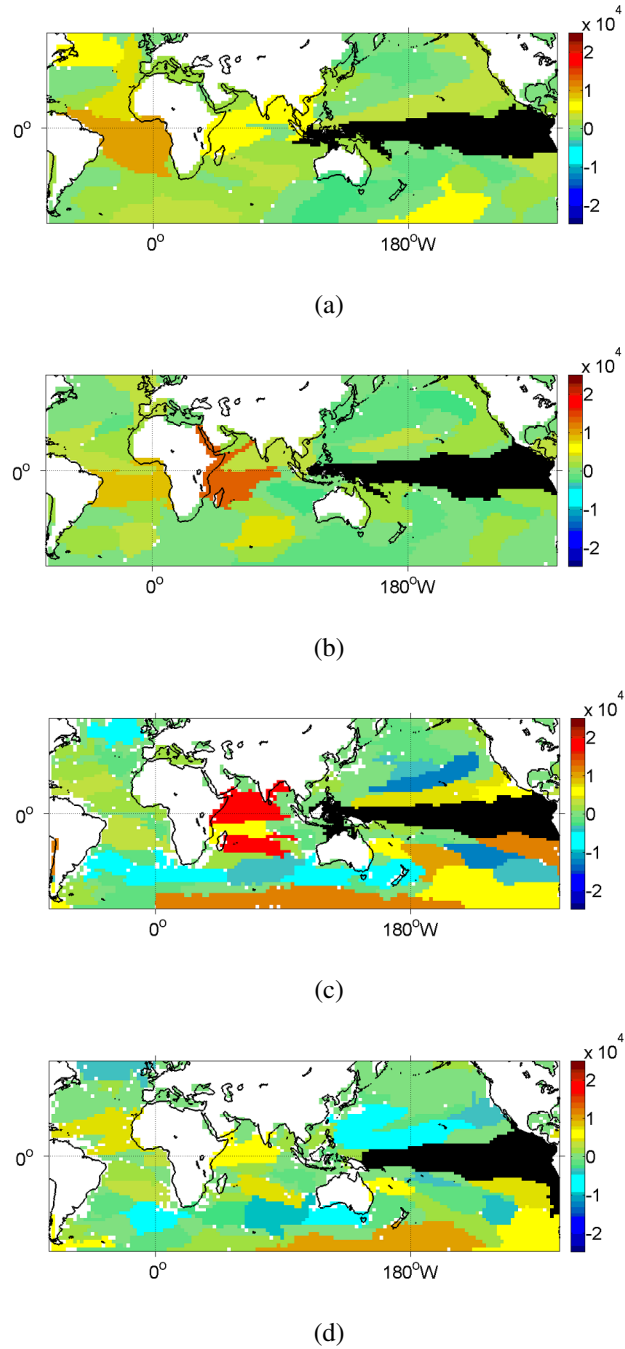


Figure 27: Link maps from the ENSO-like area in the (a-b) GISS-E2H and (c-d) HadCM3 models

0.45 for the two HadCM3 integrations. GISS again outperforms HadCM3 due to the better representation of the shape of most areas.

As already mentioned, the relative distance and adjusted Rand index metrics, while

Table 2: D_{sd} and ARI from HadISST (1979-2005) to reanalyses, GISS-E2H and HadCM3, and corresponding noise-to-signal ratios γ

Data set	D_{sd}	γ	ARI	γ
HadISST 1950-1976	0.13	0.10	0.55	0.40
ERSST-V3	0.16	0.15	0.54	0.45
NCEP	0.29	0.35	0.59	0.35
GISS run 1	0.29	0.35	0.46	0.60
GISS run 2	0.37	0.45	0.48	0.55
HadCM3 run 1	0.56	0.70	0.43	0.70
HadCM3 run 2	0.35	0.40	0.45	0.60

alone unable to quantify all the differences and similarity between networks, can be used successfully together to rank several networks with respect to a common reference. Two networks are similar if both ARI is large and D_{sd} is small, where the first constrain, given the analysis above, can be translated into $\text{ARI} \geq 0.5$ and the second into $D_{sd} \leq 0.25$. If any of these two conditions is not met, an analysis of the other metrics introduced can provide useful information on the topological differences between the data sets under consideration.

3.7 Discussion and Conclusions

We developed a novel method to analyze climate variables using complex network analysis. The nodes of the network, or areas, are formed by clusters of grid cells that are highly homogeneous to the underlying climate variable. These areas can often be mapped into well known patterns of climate variability.

The network inference algorithm relies on a single parameter τ that determines the degree of homogeneity between cells in an area. The requirement of only one parameter, combined with the fact that no link pruning in the underlying cell-level network is imposed, adds robustness to a network's structure and makes the comparison of different networks more reliable.

The constructed climate networks are complete weighted graphs. In effect, our network framework allows for investigating and visualizing the relative strength of node interactions, which can be associated with teleconnection patterns. The inferred networks are robust under random perturbations when adding noise to the anomaly time series of the climate variable under investigation, to small changes in the selection of τ , to the choice of the correlation metric used in the inference algorithm, and to the spatial resolution of the input field.

In this chapter we constructed networks for a suite of SST and precipitation data sets, and we analyzed them with a set of weighted metrics such as link maps, area strength and *s-core* decomposition. Link maps enable us to visualize all statistical relationships between areas, while strength maps highlight the relative importance of those relationships, identifying major climate patterns. The *s-core* decomposition, on the other hand, identifies the backbone structure of a network, clustering areas into layers of increasing significance. Finally, we quantified the degree of similarity between different networks using the Adjusted Rand Index metric and a newly introduced "distance metric", based on the area strength distribution.

After analyzing three SST reanalyses and two precipitation data sets, we investigated the network structure of two CMIP5 outputs, GISS-E2H and HadCM3, focusing on SST anomalies. We visualized model biases in the underlying network topology and in the spatial expression of patterns, and we quantified the distance between model outputs and reanalyses. We found significant differences between model and observational data sets in the shape and relative strength of areas. The most striking biases common to both models are the excessive longitudinal extension of the area corresponding to ENSO, and the inability to represent the horse-shoe pattern in the western tropical Pacific. Links are generally weaker than observed in the GISS-E2H model, but the relative strength, shape and size of the main areas are in reasonable agreement with the reanalyses. The HadCM3 network,

on the other hand, is closer to observations in the absolute strength of its areas, but the areas are too numerous in the tropics and unrealistically strong nodes are found in the South Pacific. In the near future, we aim at providing a comprehensive comparison of CMIP5 outputs to the climate community by extending our analysis to a much larger number of models.

In this work we limited our analysis to linear and zero-lag correlations. The methodology presented, however, could be generalized to include the analysis of nonlinear phenomena and non-instantaneous links, by introducing nonlinear correlation metrics, such as mutual information or the maximal information coefficient [122], and time-lags. Additionally, the set of metrics proposed can be enhanced to capture more complex relationships in the underlying network.

3.8 Selection of threshold τ

The threshold τ is the only parameter of the proposed network construction method. It represents the *minimum average pair-wise correlation between cells of the same area*, as shown in Eq.2. Intuitively, τ controls the minimum degree of homogeneity that the climate field should have within each area. The higher the threshold, the higher the required homogeneity, and therefore the smaller the identified areas.

Throughout this chapter, we select τ based on the following heuristic. First, we apply the one-sided t-test for Pearson correlations at level α and with $T - 2$ degrees of freedom (recall that T is the length of the anomaly time series) to calculate the minimum correlation value r_α that is significant at that level [126]. For example, with $\alpha=1\%$ and $T=81$ (corresponding to 27 years of SST monthly DJF averages), we get $r_\alpha=0.34$.

Instead of pruning any correlations $r(x_i, x_j)$ that are below r_α , we estimate the expected value of only those correlations that are larger than r_α ,

$$\bar{r}_\alpha \triangleq E[r(x_i, x_j), r(x_i, x_j) > r_\alpha] \quad (\text{A1})$$

For a set of k randomly chosen cells that have statistical significant correlations (at level

α) between them, \bar{r}_α is approximately equal, for large k , to their average pair-wise correlation. A climate area, however, is not a set of randomly chosen cells, but a geographically connected region. So, we require that the average pair-wise correlation of cells that belong to the same area should be higher than \bar{r}_α , i.e.,

$$\tau = \bar{r}_\alpha \tag{A2}$$

Note that τ is independent of the size of an area, but it depends on both α and on the distribution of pair-wise correlations $r(x_i, x_j)$.

3.9 Pseudocode of area identification algorithm

Below we present the pseudocode for the area identification algorithm used in this chapter.

```

function PART-1
  Mark all cells as available
   $k \leftarrow 0$ 
   $V \leftarrow \emptyset$ 
  while true do
    Identify the two available and neighboring cells  $(i, j)$  with the maximum correlation
    if  $r(x_i, x_j) < \tau$  then
      exit ▷ No additional areas can be identified
    else
      Area  $A_k \leftarrow i, j$ 
       $i, j \leftarrow \text{unavailable}$ 
       $V \leftarrow \text{EXPAND}(A_k)$ 
       $k = k + 1$ 
    end if
  end while
end function

function EXPAND(Area  $A_k$ )
  Construct set  $Nei(A_k)$ : all available neighboring cells to area  $A_k$ 
  while true do
    if  $Nei(A_k) = \emptyset$  then
      return  $A_k$ 
    else
       $m = \arg \max_{m \in Nei(A_k)} \hat{r}(x_m, A_k)$ , with  $\hat{r}(x_m, A_k) = \sum_{i \in A_k} r(x_m, x_i) / |A_k|$  ▷
      Identify the cell  $m$  in  $Nei(A_k)$  that has maximum average correlation with existing cells in  $A_k$ .
      if  $\hat{r}(x_m, A_k) > \tau$  then
         $A_k \leftarrow m$ 
         $m \leftarrow \text{unavailable}$ 
        Include available neighbors of  $m$  in  $Nei(A_k)$ 
      else
        return  $A_k$ 
      end if
    end if
  end while
end function

```

```

function PART-2(Areas  $V = \{A_1, \dots, A_{|V|}\}$ )
  Mark all areas  $A_i \in V$  as available
  while true do
     $A_k = \arg \max_{A_i \in V} |A_i|$        $\triangleright$  Identify the largest available area  $A_k \in V$  in terms of
    number of cells.
    if  $A_k = \emptyset$  then
      exit                                 $\triangleright$  No additional available areas.
    else
      Construct set  $Nei(A_k)$ : all geographically adjacent areas to  $A_k$ 
      if  $Nei(A_k) = \emptyset$  then
         $A_k \leftarrow$  unavailable
      else
        Identify area  $A_j \in Nei(A_k)$  such that average correlation of all cells in  $A_k \cup A_j$ 
        is maximum
        if  $\hat{r}(A_j, A_k) > \tau$  then
          Remove  $A_j$  from  $V$ 
           $A_k = A_k \cup A_j$ 
        else
          Mark  $A_k$  as unavailable
        end if
      end if
    end if
  end while
end function

```

Chapter IV

ENSO IN CMIP5 SIMULATIONS: NETWORK CONNECTIVITY FROM THE RECENT PAST TO THE TWENTY-THIRD CENTURY

4.1 Introduction

Understanding how major modes of natural variability will respond to gradual mean state changes associated with anthropogenic warming is crucial to climate science [43]. Coupled general circulation models (CGCMs) are the most powerful and widely used tool to address this problem, and therefore there is increasing interest in new approaches to evaluate systematically CGCM performances and their sensitivities to increased greenhouse gas (GHG) emissions. Here we present the first extensive application of a new methodology, built upon complex network analysis, to assess model performances in reproducing the recent past and their topological changes in future projections under varying GHG forcing.

In recent years complex network analysis [7] has been widely applied to the investigation of complex dynamical systems, ranging from the Internet and its evolution [5] to the human connectome [28]. Many of the complex systems studied by network analysis are embedded in space [16] while their elements interact with each other forming complex functional relationships. Climate is another complex system that can be represented as a spatial network. Climate networks were first introduced by Tsonis and Roebber [153], who applied ideas from graph theory to study the behavior of global geopotential height fields. Since then, climate network analysis has contributed to the discovery of new dynamical transitions and teleconnections in the climate system [94, 154, 173, 77], to the investigation of the monsoon [106, 25] and to the prediction of El Niño episodes [105]. Network approaches have been used to identify high-energy oceanic flows representing the backbone of the climate system [52], to evaluate the collective behavior of different climate

variables [51] (see also Section 3.6.4) and structural changes as climate evolves through time [20, 142, 144, 119]. Attempts to investigate causal dependencies have been made in [84] and [57].

Recently climate networks have been employed also to evaluate and compare climate models. A community detection algorithm was used to rank the performance of several CGCMs [145], and complex network analysis was adopted to evaluate the Statistical Analogue Resampling Scheme (STARS) model against a dynamical model (COSMO-CLM) in representing the climate of South America [61].

Here we analyze the network properties of model outputs from the Coupled Model Intercomparison Project - Phase 5 (CMIP5) spanning the 1956-2100 or 1956-2300 intervals using the network methodology proposed in Section 3.3 of this thesis. First we identify areas - i.e. geographically connected regions homogeneous to the underlying climate variable - that represent the nodes of the network, roughly corresponding to major climate modes. Then we visualize, validate and compare those areas and their links or connections. Links represent non-local dependencies between different areas. Therefore, in contrast to more commonly used community detection techniques, our method decouples the identification of climate nodes from the connections that those have with each other.

The methodology adopted yields several desirable properties compared to more traditional time series analysis [10, 1, 11, 43, 49, 62, 73]. It allows evaluating model performances at both local and global scales, uncovers relations in the climate system that are not fully captured by traditional methodologies, explains known climate phenomena in terms of the underlying network's structure and metrics, and is not locked into a particular set of climate indices from the outset. Its scope is similar to empirical orthogonal functions (EOF) in that it identifies the major modes of climate variability for a given variable. In contrast with EOFs, however, our method does not impose any orthogonality constraints and does not mask patterns (or climate modes) of weaker variance. Regional or global changes can be quantified in terms of addition or removal of areas, fluctuations in the weight of existing

links, or variations in the relative significance of different areas, providing sensitivity information. The proposed framework is fast and scalable, and has been developed to ensure robust comparisons. Furthermore, it allows estimates of model trajectories over time and of intra ensemble variability. The last can be objectively compared to contributions from different forcings or mean states.

In this work we focus on global quantities, and analyze time intervals of fifty years. The dominant mode of variability at those time and spatial scales is the El Niño Southern Oscillation (ENSO). First we assess how CMIP5 models represent the network topology of ENSO and its teleconnections in sea surface temperature and precipitation comparing them to various reanalysis. Then we focus on model projections and on the stability of ENSO and its links in the near and far future.

4.2 Climate Network Inference

The network inference is a three-step process. First we construct a “cell-level network”; second we apply a clustering algorithm to identify the nodes or areas; third we compute weighted links between areas to quantify their connections. All networks are inferred from monthly averages of detrended seasonal anomalies of sea surface temperature (SST) and precipitation but the procedure can be applied to any variable of interest. Trends are calculated with the Theil-Sen estimator [6] to reduce sensitivity to outliers and at least partially account for the ENSO variability [135]. All datasets are interpolated to a minimum common resolution ($2^\circ lat \times 2.5^\circ lon$ for SST and $2.5^\circ lat \times 2.5^\circ lon$ for precipitation) and only the range $[60^\circ N - 60^\circ S]$ is considered due to the large differences in reconstructions, reanalyses and models at higher latitudes. We focus on boreal winter (December to February, DJF), when ENSO is strongest. Calculations have been repeated for summer (June to August) confirming all major outcomes.

The cell-level network is constructed computing the Pearson cross-correlation $r(x_i(t), x_j(t))$ between the anomaly time-series $x_i(t), x_j(t)$ for all cell pairs i, j . All pair correlations are

retained and the resulting cell-level network is a complete weighted graph (i.e. a link exists between all pairs of grid cells). This characteristic differentiates this method from most prior work on climate network analysis where a threshold to prune non-significant correlations is applied [53, 174, 143, 155], and guarantees robust comparisons between networks constructed on different datasets. The cell-level network is input to the clustering algorithm and relies on a single parameter τ controlling the homogeneity of areas to the underlying climate variable. Formally, an area A_k is a geographically connected cluster of two or more cells satisfying

$$\frac{\sum_{i \neq j} r(x_i(t), x_j(t))}{|A_k|(|A_k| - 1)} > \tau, \quad (8)$$

where $|A_k|$ denotes the number of cells in the area. τ represents the minimum average pair-wise correlation between cells of an area at a given significance level α (here $\alpha = 1\%$) and is determined following the heuristic presented in Section 3.8. τ depends on α and on the distribution of pair-wise correlations $r(x_i(t), x_j(t))$ in any given dataset.

The clustering algorithm aims also to minimize the number of areas identified; the problem is NP-Complete, thus the algorithm relies on greedy heuristics. It identifies areas iteratively by selecting the pair of geographically connected grid cells with the maximal $r(x_i(t), x_j(t))$; An area is further expanded by adding the adjacent grid cell that maximizes the average cross-correlation to the existing cells in the area. The area expansion stops either if Eq.8 is violated or when all neighboring grid cells belong to other areas¹. Since the algorithm relies on greedy heuristics, the solution is suboptimal and areas with at least one pair of geographically adjacent cells, and whose union satisfies Eq. 8, are further merged together. The methodology ensures the robustness of the area-level structure for a wide range of significance levels, as extensively tested (see Section 4.6 and Section 3.5.3).

Finally, links are computed from the area cumulative anomalies. For a given area A_k , the cumulative anomaly is equal to $X_k(t) = \sum_{i \in A_k} x_i(t) \cos(\phi_i)$, with ϕ_i being the latitude

¹At the end of the area identification some grid cells may not belong to any area (if they violate the τ criterion for each candidate area). Such grid cells are shown in white in all maps.

of cell i (the anomaly time series of any given cell i are therefore weighted by the cell size). The weighted link $w(A_k, A_m)$ between two areas A_k and A_m is equal to the covariance between the corresponding cumulative anomalies. Links can be positive or negative, and are computed for all pairs of areas to obtain a complete weighted graph. Link maps allow the visualization of the (weighted) connections between any given area and all others in the network. Areas are also characterized by their weighted degree or strength, defined as the sum of the absolute link weights $W(A_k) = \sum_{m \neq k}^V |w(A_k, A_m)|$, where V is the set of the areas $A_1 \dots A_V$ inferred. Strongest areas correspond to major modes of climate variability.

Similarities and differences between two networks N and N' , each of size n grid cells, are quantified by two metrics, the Adjusted Rand Index (ARI) and a newly defined network distance D . The ARI measures the spatial likeness of the areas in two networks [86, 140]. Any pair of cells that belong to the same area in N and N' , or that belong to different areas in both networks, contributes positively to the ARI; conversely, any pair of cells that belong to the same area in one partition but to different areas in the other, contributes negatively. The ARI ranges between 0 and 1, with 1 denoting perfect similarity, and ensures that the distance between two random partitions is zero. The ARI, however, does not consider cell anomalies and (actual) cell size.

To capture similarities or differences at the network level (i.e. in terms of link weights and area strengths) we also define a distance D between two networks. For the calculation of D we assign each grid cell a weight that is equal to the strength of the area the cell belongs to. The distance D between two networks N and N' is then defined as

$$D(N, N') = \frac{\sum_{i=1}^n |W_N(i) - W_{N'}(i)|}{\sum_{i=1}^n |W_{\hat{N}}(i) - W_{\hat{N}'}(i)|} . \quad (9)$$

n is the number of grid cells and it includes cells that do not belong to any area. $W_N(i)$ is the weight assigned to grid cell i in network N ; similarly for $W_{N'}(i)$. The network \hat{N} is a randomized instance of N in which the cells of the latter have been randomly permuted in the underlying grid, keeping the original weight that was assigned to them in N . The numerator of D increases whenever a cell belongs to different areas, or, if two areas are

identical, whenever they have different strengths. The denominator of D is expected to be higher than the numerator because it is very unlikely that the same grid cell of the two randomized networks belongs to the same area. In the pathological case that two networks differ significantly in terms of their areas but all grid cells have roughly the same weight, the distance D will still be high (close to one, given that the numerator and denominator will be approximately equal). It is noted that D is different than the distance metric that was introduced in Section 3.4; the metric of Eq. 9 considers not only the strength of each area but also its spatial extent, while the distance metric of Section 3.4 only considers the area strength distribution.

The joint consideration of both ARI and D offers more information than any of the two metrics alone. Specifically, ARI focuses on the spatial extent of each area (the set of cells that belong to an area) but it ignores the area strengths. The distance D depends both on the spatial extent and the strength of each area but it does not separate the two. So, for instance, when two pairs of networks both have $D \sim 0$ but one of them has higher ARI, we can conclude that the latter are more similar compared to the other pair, mostly due to the spatial extent of their areas.

Finally, given two networks N and N' and their respective $D(N, N')$ and $\text{ARI}(N, N')$, it is possible to map both metrics to the amount of white Gaussian noise (WGN) that added to the original climate field will produce a network N'' such that $D(N, N') \approx D(N, N'')$ and $\text{ARI}(N, N') \approx \text{ARI}(N, N'')$. Specifically, the anomaly time series $x(t)$ of the original climate field can be perturbed by adding WGN γ -times the variance of $x(t)$. γ therefore quantifies the noise-to-signal ratio between N and N' .

Several different approaches have been proposed in the literature to represent the Earth's climate as a network. A common element in most of them is that the network nodes are grid cells and edge pruning is performed to remove non-significant pairwise correlations. Our methodology differs substantially, and in the following, we contrast it with the two most relevant climate network methods developed to assess climate model outputs [61, 145]. In

[61] the authors evaluate the performance of two regional models representing the South American climate. Their method represents the climate network as a binary graph. Nodes correspond to grid cells, weighted proportionally to their geographical size. Non-significant links are removed by enforcing a fixed graph density and only positive correlations are considered. In [145] the authors evaluate the performance of an ensemble of CMIP3 models. The climate network is again represented as a binary graph. In contrast to [61], both positive and negative correlations between nodes are taken into account. Network nodes are unweighted and non-significant edges between them are removed using a fixed threshold approach. The climate network is then used as input to a community detection algorithm. A community is a subset of nodes that are densely interconnected relative to their connections with the rest of the network. The identified communities are groups of grid cells forming, possibly disjoint, geographical regions. Model differences are captured using the ARI metric, measuring the spatial similarity between the identified communities in each network. Summarizing, in [61] models are evaluated based on their actual network structure, while in [145] model outputs are evaluated based on their community structure.

Instead, the proposed methodology compares climate models based both on network structure (distance metric) and on the spatial representation (ARI metric) of different climate modes of variability. Furthermore, the combination of the ARI and D metrics allows also to quantify intra-ensemble variability, while modeling the climate network as a weighted graph enable to evaluate the magnitude and relative importance of specific teleconnections. By considering both positive and negative link weights, different functional relationships between the elements of the climate system are considered. Similarly to community detection, grid cells are clustered into areas and this reduces the dimensionality of the problem. However, communities may consist of geographically disjoint areas, and so they will not show explicitly the teleconnections between these regions. In contrast, we decouple the identification of areas from the connections they have with each other.

4.3 Results

4.3.1 CMIP5 Models and Observational Datasets

The network analysis is performed on realizations from twelve models of the CMIP5 catalog (Table 3), chosen among those with ensembles of at least three members in the historical period, and with one member or more continuing to 2100, and possibly to 2300, under the scenario with the highest Representative and Extended Concentration Pathways (RCP8.5 and ECP8.5) relative to preindustrial levels [109]. The projections are forced with emissions such that the radiative forcing induced by GHGs reaches -8.5 Wm^{-2} in 2100 [124]. This choice of scenario is dictated by the larger availability of modeling centers extending their integrations to 2300. To evaluate the realism of CMIP5 CGCMs in simulating the recent past, we consider historical ensembles over the period 1956 - 2005 [149], and we contrast SST and precipitation model networks with the ones from the Hadley Center SST reconstruction over the same period (HadISST) [121], and from the European Centre for Medium-range Weather Forecasts Re-Analysis (ERA40+Interim). ERA40+Interim combines ERA-40 [159], available from 1958, with ERA-Interim [46] after 1979. Furthermore, networks constructed from the Extended Reconstructed Sea Surface Temperature version 3 (ERSST-V3) [134], surface temperatures provided by National Centers for Environmental Prediction (NCEP) [93], and two SST realizations of the Simple Ocean Data Assimilation reanalysis (SODA version 2.1.6 available from 1958 to 2005, and 2.2.8, covering the 1956-2005 interval) [31], are compared to HadISST to quantify the range of uncertainties and spread in the SST observational proxies. For precipitation, networks from the NCEP reanalysis, the CPC Merged Analysis of Precipitation (CMAP) [172] and ERA-Interim, the last two available from 1979, are also compared to ERA40+Interim. We verified that NCEP rainfall networks over the 1958-2005 or 1956-2005 periods are indistinguishable. Networks are then computed for the model future projections, and for all integrations, past and future networks are compared to quantify projected changes in climate modes (areas) and their connections (links). Similarities and differences between HadISST and CMIP5

historical SST networks, ERA and modeled precipitation networks, and between historical and projected networks for the same model member, are summarized using the Adjusted Rand Index and the distance D .

Climate networks are constructed using detrended time series of SST and precipitation. One question we wish to answer is if the uncertainty in the representation of major tropical teleconnections in CMIP5 models results in greater or lesser regional impacts than the uncertainty in temperature and rainfall trends. Therefore a brief comparison of observed and modeled trend for the historical period, and a description of future trends during the projected intervals is added to each subsection, prior of the network analysis.

4.3.2 The Historical Experiments: 1956-2005

Historical global mean trends in winter are summarized in Table 3 for models and observational proxies. Several models overestimate the observed SST trend over the historical period, due to their inability to simulate the 'pause' or 'hiatus' observed since 1998 [70]. In the majority of integrations SSTs are characterized by cooling (or lesser warming) south of $50^{\circ}S$, that results from heat uptake by the deep ocean [97], and by the greatest warming over the Atlantic and Indian Oceans between $40^{\circ}S$ and $50^{\circ}S$, in agreement with the observational proxies (Fig. 30, left panels). Most models warm above the global mean in the Equatorial Pacific and show negative trend anomalies in the East China Sea and along the coasts of Japan. Conversely, the observational proxies display a cooling trend along the equatorial Pacific, and the most intense cooling in the central North Pacific [99], and in the subpolar gyre in the North Atlantic. Global mean precipitation trends are extremely small for models, and uncertain in the reanalyses (Fig. 30, right panels). At a regional level, however, NCEP and ERA (and CMAP over the available period) have slopes much steeper than any CMIP5 output, and a complex spatial patchiness that varies greatly with the period considered, indicating large interannual fluctuations not represented in the CGCMs. In the tropics, all models but CSIRO and MIROC5 underestimate the local trends by two- or

threefold. In the extratropics none of the runs captures the observed variability, underestimating it by five times or more. 70% of models show an increase in rainfall in the tropical Pacific centered around $5^{\circ}S$, in partial agreement with ERA.

Table 3: List of models analyzed and global mean trends in sea surface temperature and rainfall over 1956-2005 and 2051-2100. The number of ensemble members considered during the historical period (1956-2005) is indicated for each model. In parenthesis the number of members with projections to 2100 under the RCP8.5 scenario. X indicates that the model has one member continuing to 2300. Boreal winter (December to February) global mean trends are averaged over all ensemble members (\pm denotes the maximum deviation between ensemble members)

Model	Ensemble #	1956-2005 SST $C^{\circ}/\text{year} \times 10^{-2}$	1956-2005 PREC (mm/day)/yr $\times 10^{-4}$	2051-2100 SST $C^{\circ}/\text{year} \times 10^{-2}$	2051-2100 PREC (mm/day)/yr $\times 10^{-4}$
BCC-CSM1.1	3(1) X	1.2 ± 0.3	9.5 ± 2.7	3.2	26.7
CanESM2	4(4)	1.2 ± 0.2	8.4 ± 2.2	4.0 ± 0.9	21.0 ± 2.8
CCSM4	4(4) X	1.4 ± 0.1	9.6 ± 2.6	3.5 ± 0.2	23.0 ± 3.4
CNRM-CM5	4(4) X	0.7 ± 0.4	1.1 ± 2.4	3.5 ± 0.1	25.0 ± 3.8
CSIRO-Mk3.6.0	4(4)	1.2 ± 0.1	-2.0 ± 2.6	4.1 ± 0.1	34.0 ± 3.7
GFDL CM3	4(1)	0.8 ± 0.2	2.2 ± 2.6	4.2	31.0
GISS-E2-H	4(4) X	0.6 ± 0.3	2.6 ± 3.7	2.4	14.0
HadGEM-ES	4(4) X	0.5 ± 0.3	1.4 ± 1.6	4.0 ± 0.3	21.0 ± 2.8
IPSL-CM5a-LR	4(4) X	1.4 ± 0.1	14.0 ± 1.6	4.4 ± 1.3	39.0 ± 7.5
MIROC5	4(3)	0.7 ± 0.1	2.6 ± 2.1	2.9 ± 0.1	17.0 ± 1.5
MPI-ESM-LR	3(3) X	1.0 ± 0.1	11.0 ± 1.0	3.3 ± 0.1	25.0 ± 3.2
MRI-CGM3	4(1)	0.6 ± 0.2	-0.1 ± 3.3	3.1	31.0
REANALYSIS		1956-2005 HadISST $C^{\circ}/\text{year} \times 10^{-2}$	1958-2005 ERA (mm/day)/yr $\times 10^{-4}$	1958-2005 NCEP (mm/day)/yr $\times 10^{-4}$	
		0.7	53.1	-3.1	

Comparing the observational proxy networks in terms of their global metrics, the NCEP surface temperature reanalysis is further apart from HadISST than any other observational dataset, with $\gamma > 1$ (Fig. 31a). This was expected considering that we are comparing SST with surface air temperature (masked over the ocean), and can be used as benchmark for the model comparisons. In particular, the NCEP reanalysis overestimates the strength of the areas covering the tropical Indian Ocean, and misrepresents the so-called horse-shoe pattern in the Pacific (see strength maps in Section 4.5, Fig. 42).

Between the models, seven have at least one realization contained within the uncertainty cloud of the reanalyses, as MIROC5, shown in Fig. 32b, that displays a network slightly

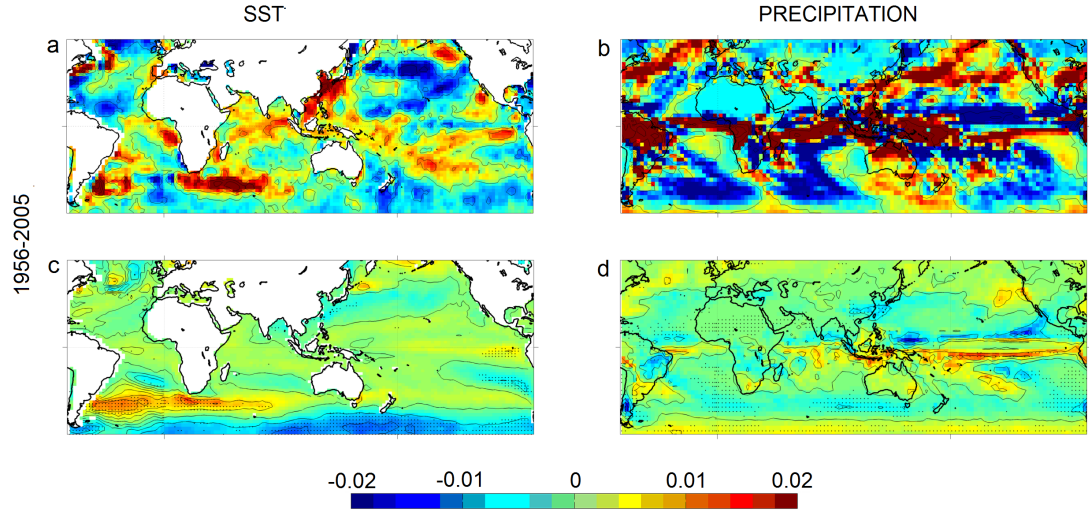


Figure 30: Trend anomaly maps for boreal winter in the recent past and near future. Anomalies are computed by removing the global mean trend calculated over the months of December to February and indicated in Table 3 from each grid cell. + and • indicate agreement in more than 90% and 70% of models in the sign of the trend anomaly slope. (a) HadISST. (b) ERA40+Interim. (c) Sea surface temperature (SST) averaged across models in the historical period (1956-2005). (d) As in (c) but for rainfall. The units are C°/year for SST and $(\text{mm}/\text{day})/\text{year}$ for precipitation

stronger but overall very similar to the observed. Of the remaining BCC, GISS-E2H and MRI underperform in both metrics due to an underestimation of size and strength of most areas (see Fig. 42d for MRI and Fig. 42 for one sample map from each modeling center) and very weak connectivity between nodes (e.g. Fig. 33d and Fig. 43). Furthermore, the area corresponding to ENSO develops too narrowly around the Equator, extends into the Warm Pool region, and has low strength, which is directly associated to a very low ENSO variance. The extension of the ENSO node into the west Pacific is common to HadGEM2, but strength and connectivity of major areas compare well to observations. GFDL CM3, shown in Fig. 32c, and IPSL display a strong, broad area in the Southern Ocean (SO) south of $45^{\circ}S$ extending from the Atlantic to the Pacific. BCC, CanESM2 and CNRM display an analogous node, but of lesser strength. In IPSL the SO area has comparable or greater strength than ENSO. The correlation between the SO node and ENSO is zero or moderately positive in IPSL, and generally very high and positive in GFDL CM3 (Fig. 33c

and in Figs. 42 and 43 in Section 4.5).

The spread in ARI and D for members of the same ensemble is indicative of the model intrinsic variability. In general, large intra-model differences are noticeable for CGCMs further apart from HadISST and are related to the strength of the areas. D is strongly affected by the connections that the ENSO-related node reproduces: Models with weak ENSO areas (e.g. BCC, GISS-E2H, MRI) or for which ENSO is not the strongest mode of variability (IPSL), are subject to greater spreads in their distance. One member may yield nodes that are too weak, while representing correctly their relative strengths and links, and another member may develop implausible relations between areas other than ENSO. Coupled models capable of reproducing well the strength of the ENSO node, and for which this node is dominant in the network, cluster their members closely together, even more so that different observational proxies.

For precipitation, the two strongest areas observed in the tropical Pacific correspond to ENSO and the Warm Pool (bottom rows in Fig. 32, and Fig. 44). In all reanalyses and CMAP the node associated to ENSO extends along the equator from about $180^\circ W$ to the coast of the American continent. The spread in D between different rainfall reanalyses is far larger than for SSTs, with NCEP displaying the least agreement with ERA40+Interim (Fig. 31b). Additionally, the ARI is always smaller than 0.5, indicating profound differences in the node shapes and distributions also between datasets representing the observational truth. Precipitation is by nature an intermittent field in space and time; the inferred networks have a much greater number of nodes than their SST counterpart, reducing the chances of spatial and strength likeness. The quantification of the noise to signal ratio accounts for this inherent difference between SST and rainfall metrics. The comparison of Fig. 31a and 31b reveals that models with SST networks characterized by very large D and small ARI perform poorly also in representing rainfall. BCC, GISS-E2H, IPSL and MRI underestimate, once more, the strength of most tropical areas (see for example Fig. 32h). Additionally, the strongest node in BCC and IPSL occupies the center of the Pacific Ocean

and does not penetrate eastward, in MRI extends from $180^{\circ}W$ to the west, reaching New Guinea, and in two of the GISS-E2H members fills the whole equatorial Pacific. A reliable representation of SST variability, however, does not guarantee a realistic simulation of precipitation distribution and interannual modulation. CSIRO and MPI, in particular, are penalized in the metrics due to the shift of their rainfall ENSO-related area westward, over the center of the Pacific basin; furthermore, CSIRO underestimates the size of major nodes. CNRM outperforms all other models (and partially NCEP) with both smallest D and largest ARI, followed by MIROC5 (Fig. 32f). They both reproduce well the patterns associated to ENSO and the Warm Pool in the equatorial Pacific, in terms of shape, and more so strength, and are capable of simulating major connections between nodes (Fig. 45). In all other CGCMs the area corresponding to the Warm Pool is absent, as in HadGEM2, or shifted to the west into the Indian Ocean, as in MPI. Finally, different ensemble members appear clustered together more tightly than the reanalyses and CMAP, independently of their ability to represent the observations. As a result, the wide range of strengths found in SSTs for models with a weak ENSO node is not mirrored in precipitation.

4.3.3 The RCP8.5 Experiments: 2051-2100

Near future trends are projected to be analogous in patterns, but stronger in amplitude, to those found during the 1956-2005 period in both SST and rainfall (Fig. 34).

Globally CGCMs warm by $3.5 \times 10^{-2} C^{\circ}/\text{year}$ on average, and get wetter (Table 3). They agree in the main on the areas subject to above average warming (North Pacific sub-polar gyre, equatorial Pacific, Arabian Sea, the band between $40^{\circ}S$ and $50^{\circ}S$) and cooling (south of $50^{\circ}S$, the eastern side of the South Pacific gyre, the eastern side of the North Atlantic); or to more intense rainfall (north equatorial Indian Ocean, south equatorial Pacific around $5^{\circ}S$, North Pacific gyre) and weaker precipitation (regions to the north and south of the Pacific interconvergence zone).

Figure 35 summarizes the differences between SST and precipitation networks over

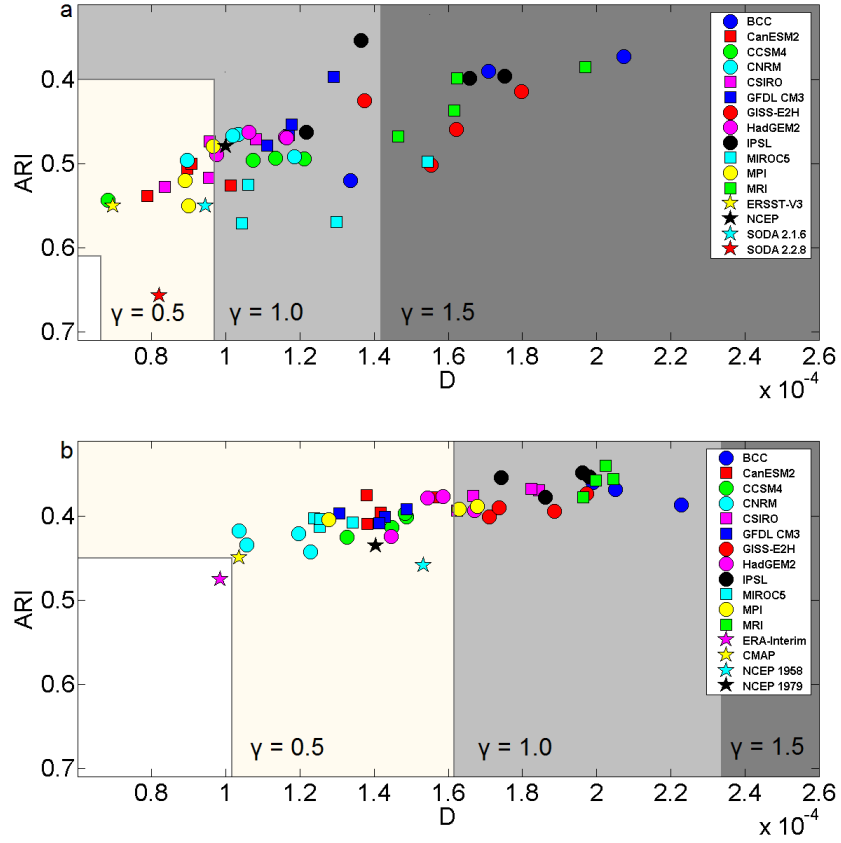


Figure 31: Metric D versus ARI for climate networks during the historical period 1956-2005. (a) Sea surface temperature; reference network HadISST. (b) Precipitation; reference network ERA40+Interim. Three levels of noise-to-signal ratios γ are also indicated

2051-2100 from their historical counterparts. Focusing on SST, all models with more than one integration available (i.e. all but BCC, MRI and GFDL CM3) have at least one member whose projected areas into the 21st century closely resemble those found in the historical period. An example from the CanESM2 model is shown in the left panels of Fig. 36 (see also Fig. 46 for a sample map for each modeling center). For those projections D and ARI from the corresponding 20th century realization are contained within the spread of the reanalyses (i.e. $D \leq 10^{-4}$ and $[ARI \geq 0.5]$); The changes in topological properties and connectivity (Fig. 47) are therefore insignificant, and the response to increased GHGs is simply the superposition of the trends, with their regional patterns, onto their historical modes of

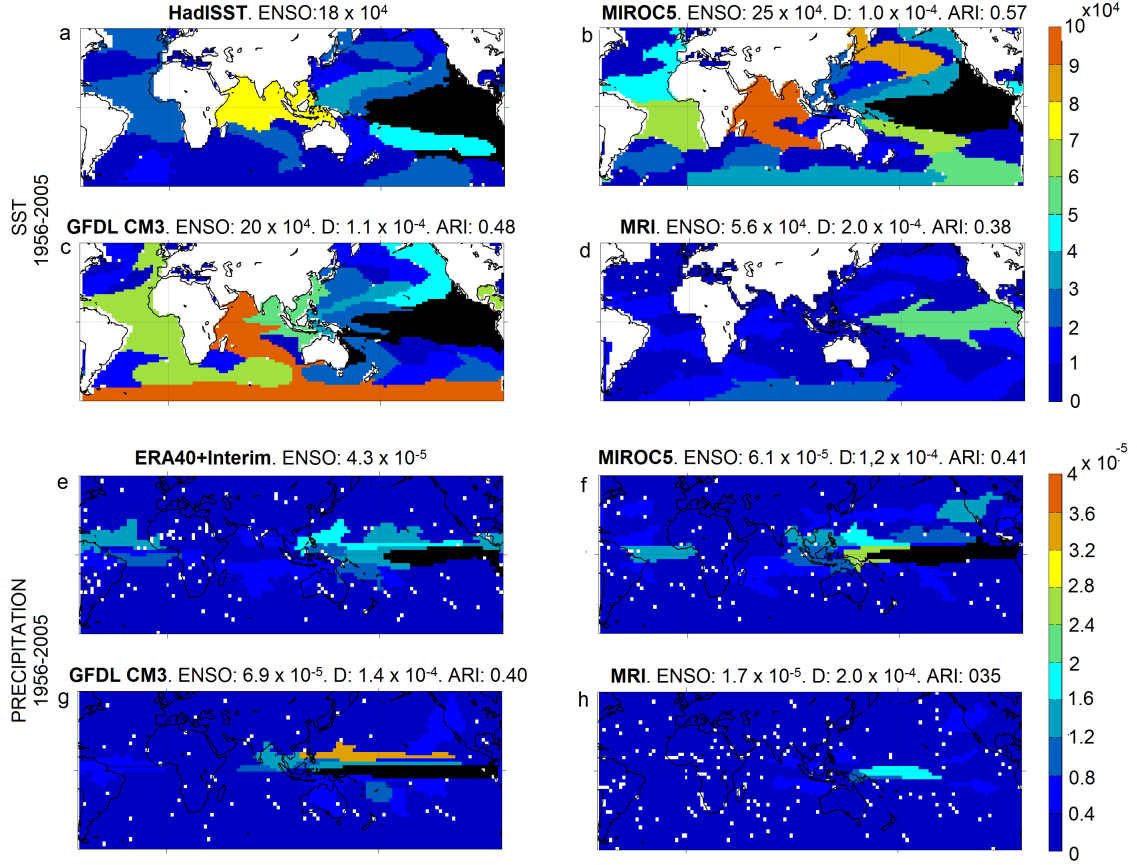


Figure 32: Strength maps of sea surface temperature for HadISST and three sample models (top rows), and of precipitation for ERA40+Interim and the same three models (bottom rows) during the historical period 1956-2005. Models shown: MIROC5, GFDL CM3 and MRI. For clarity, the strength of the ENSO-related area is saturated when exceeding the colorscale and its value is indicated at the top of each panel, together with D and ARI from HadISST or ERA40+Interim for each of the model networks

variability. Of the remaining, GFDL CM3 displays a significant change in spatial likeness due to the disappearance of the Southern Ocean node. In eight out of twelve models the members that differ in D display a decrease in strength of the ENSO area and its connectivity by a third or more of the historical value, as for the member of CanESM2 shown to the right in Fig. 36. The same eight models are characterized by a more prominent tendency for eastward propagation of positive ENSO events, associated with a weakening of the equatorial upper ocean currents, as noticed by [129]. The exceptions are MIROC5, MPI and MRI, where the ENSO node strengthens, and IPSL, where the ENSO area weakens to

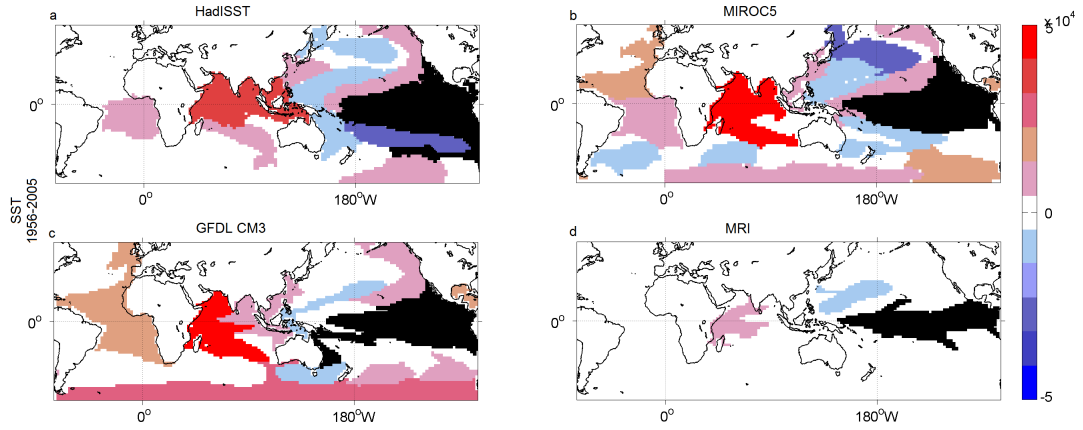


Figure 33: Sea surface temperature link maps from the ENSO-related area in black for HadISST and the three sample models during the historical period 1956-2005. Models shown: MIROC5, GFDL CM3 and MRI

a small degree, and the Southern Ocean node becomes stronger and dominant. In those four models moderate or no changes in propagation asymmetry have been found [129], but the MIROC5 version analyzed differs from ours.

Precipitation networks for the RCP8.5 scenarios do not differ from their historical counterparts more than the reanalyses and CMAP over the historical period, in both ARI and *D*. Only MRI and one member of GISS-E2H stand out due to a large increase in strength of the node associated to the ENSO anomalies in the equatorial Pacific and increased connectivity (Fig. 48 and 49). In MRI the ENSO related area is five times stronger than in the historical period, pointing to a considerable sensitivity of the model convective scheme to SST changes, and in the GISS-E2H member is almost three times stronger, achieving a value close to the reanalyses over the historical period.

4.3.4 The ECP8.5 Experiments: 2101 - 2300

For models simulating the climate system evolution under the highest of the Extended Concentration Pathways (ECP8.5) [109], the network analysis is extended to 2300. According to the ECP, the aggregated GHG emissions rise until 2100, remain constant until 2150, drop linearly to current levels by 2250 and continue as such to 2300. Correspondingly,

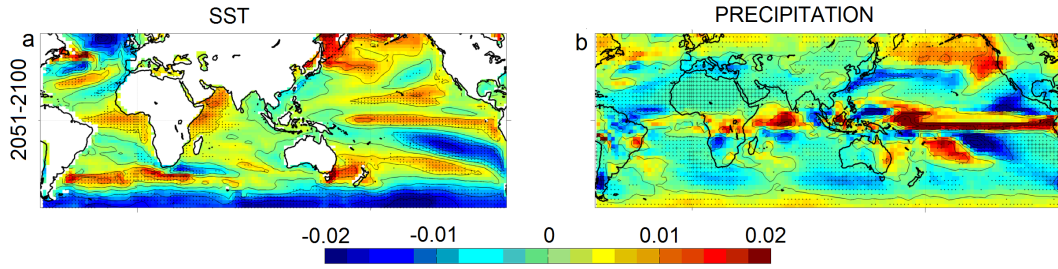


Figure 34: Trend anomaly maps for boreal winter in the second half of the 21st century. Anomalies are computed by removing the global mean trend calculated over the months of December to February and indicated in Table 3 from each grid cell. + and • indicate agreement in more than 90% and 70% of models in the sign of the trend anomaly slope. (a) SST averaged across models over 2051-2100. (b) As in (a) but for rainfall. The units are C°/year for SST and $(\text{mm}/\text{day})/\text{year}$ for precipitation

the warming trend decreases with time, especially in tropical regions (Figure 37 and Table 4). Seven of the twelve models have one member continuing to 2300. The networks are constructed on four consecutive fifty-year windows.

Figure 38 presents D and ARI for SST and precipitation, again evaluated against their historical counterpart. For clarity, the distance for the corresponding ensemble member during 2051-2100 is repeated. The SST networks for BCC, CCSM4, CNRM, GISS-E2H, and HadGEM2 depart significantly from the historical period, and they are characterized by increasingly greater distances, exceeding $\gamma > 1.5$ by 2150 or 2200. The large distances are due to a decrease in strength of the ENSO-related area and its links to half or a quarter of their original value (Section 4.5, Fig. 50 and 51). None of these models recovers ENSO and its teleconnections once emissions are reduced. In fact the ENSO area in CCSM4 first expands west into the Warm Pool region while retaining its strength and major links (2051-2150), and then weakens dramatically and suddenly after 2150 (Fig. 39a,c), in HadGEM2 loses its strength after 2200, and in GISS-E2H it is not the dominant mode of variability past 2250. A different trajectory is followed by IPSL with a reduction in strength in the tropics that culminates in 2200 and is partially recovered by 2300. Through the whole integration IPSL produces a network with a strong SO area, which at times - from 2050 to 2200

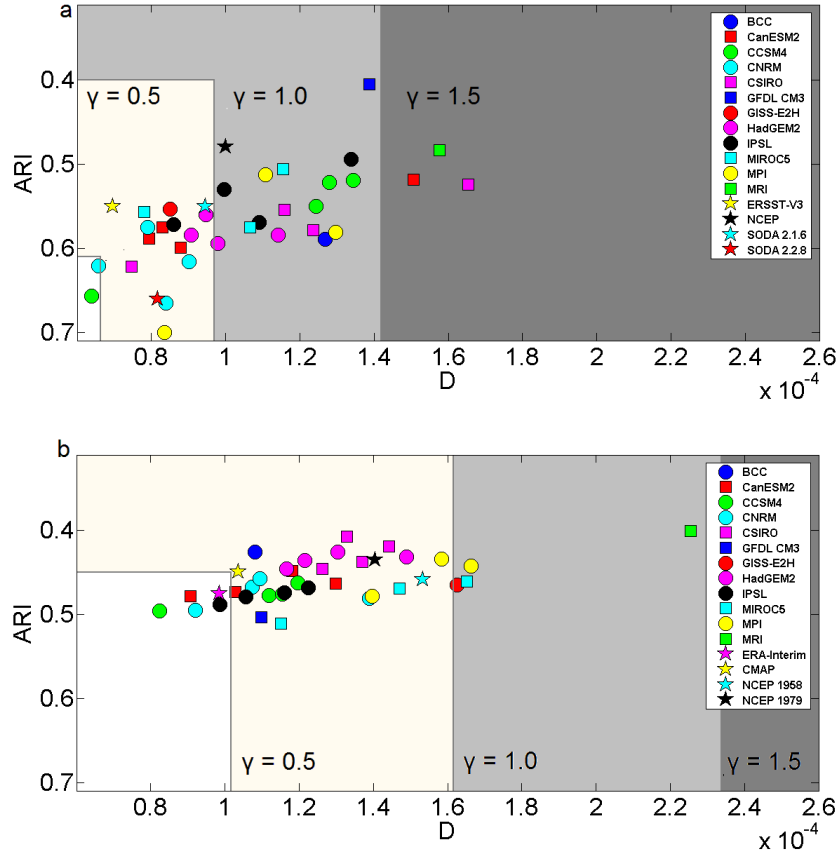


Figure 35: Metric D versus ARI for climate model networks during the period 2051-2100. (a) Sea surface temperature. (b) Precipitation. All networks are referenced to the corresponding integration over the historical period. Three levels of noise-to-signal ratios γ are also indicated. D and ARI between HadISST and other sea surface temperature proxies, and ERA40+Interim and other precipitation reanalyses are repeated to provide context

- is stronger than the ENSO node. Finally, MPI responds to the warming by strengthening the ENSO area and its links, particularly over the Indian Ocean and the tropical Atlantic in the 21st century, and oscillating between a network stronger than, or comparable to, its historical counterpart in the following periods. After 2200 the differences between historical and projected networks are negligible in strength (less or equal to differences between observational proxies) and minor in area likeness (Fig. 39b,d), with the ENSO node no longer extending into the Warm Pool. The differences in the evolution of the strength of the ENSO area are reflected in its connectivity: links from the ENSO node are dramatically

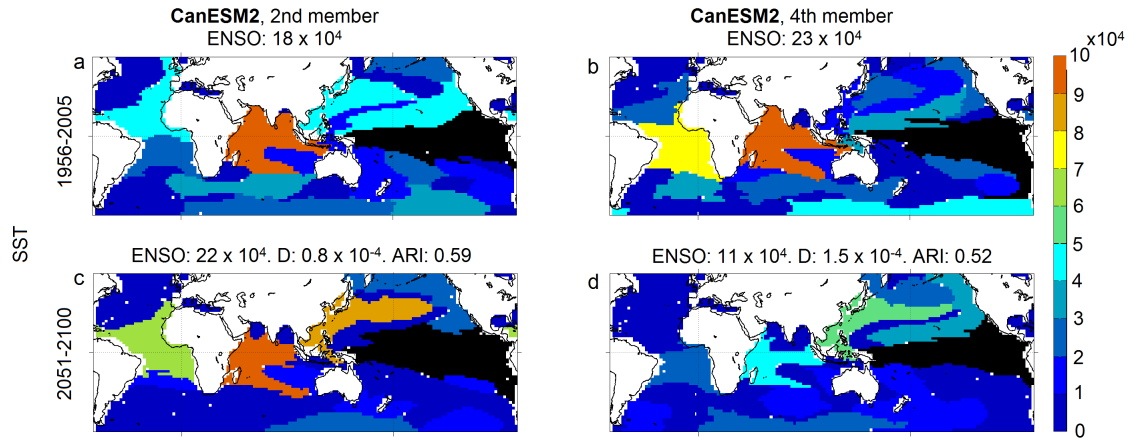


Figure 36: Sea surface temperature strength maps for two members of the CanESM2 model in the historical period (1956-2005) on top, and in the 21st century (2051-2100) at the bottom. For clarity, the strength of the ENSO-related area is saturated when exceeding the colorscale and indicated in each panel. In the future projections D and ARI from the corresponding historical member are also specified

reduced in most models, while remain comparable to the recent past in MPI (Fig. 40a,b).

The time series of the cumulative anomaly over an area quantify the evolution of its strength variance. For ENSO the variances in the historical period and during 2251-2300 in DJF are shown in Fig. 41, with the HadISST plotted as reference. All models but MPI display a systematic, gradual reduction in mean variance, ranging from -33% in GISS-E2H to -75% in CCSM4 by 2300. Large changes in ENSO variance ($\pm 50\%$) have been found also for millennial unforced simulations [48, 37], but without a preferred sign tendency, while fossil corals suggest that a weaker ENSO than today dominated the last 10,000 years [37, 30]. The ENSO variance in the historical period varies depending on the twenty-year window used for the calculation, as indicated in Fig. 41. Such variability is twice as strong as the observations in MPI and about half as observed in BCC and GISS-E2H, while is consistent with the reanalysis in the remaining four models.

In the case of precipitation, the networks for BCC, CNRM, and HadGEM2 are unaltered in the projections, except for a mild weakening of most tropical areas in the first two models (Fig. 52). GISS-E2H, after an initial strengthening, returns to conditions close

to the historical period by 2300. CCSM4 exhibits a fivefold decrease in the strength of the ENSO area over two hundred years, while the nodes covering the south Pacific convergence zone and the south equatorial Indian Ocean become stronger (Fig. 39e,g). Those areas eventually lose their connectivity with ENSO and evolve independently of it (Fig. 53). In IPSL the strength of the ENSO node fluctuates, decreasing slowly at first and regaining power in the last 50 years. The area also shifts position, translating eastward and occupying first the western and central Pacific, then the central portion, and finally developing to the east of $180^{\circ}W$ after 2250. Finally, MPI after strengthening most nodes in the 21st century by almost three folds, and shifting the strongest one eastward, maintains the new strengths and intensifies the links between ENSO and all major areas, from the Warm Pool to the Indian Ocean, and the north and south Pacific (Fig. 39f,h and Fig. 40d). By 2300 the rainfall network is much stronger and more complex than during the historical period, in spite of the SST network resembling the 20th century one.

Table 4: Projected global mean trends in sea surface temperature and rainfall from 2101 to 2300. Trends are calculated over 50-year long consecutive intervals for the models with one member extending to 2300 and for boreal winter (December to February). Precipitation trends are in parenthesis

Model	2101-2150 SST [PREC] C°/year $\times 10^{-2}$ [(mm/day)/yr] $\times 10^{-4}$	2151-2200 SST [PREC] C°/year $\times 10^{-2}$ [(mm/day)/yr] $\times 10^{-4}$	2201-2250 SST [PREC] C°/year $\times 10^{-2}$ [(mm/day)/yr] $\times 10^{-4}$	2251-23000 SST [PREC] C°/year $\times 10^{-2}$ [(mm/day)/yr] $\times 10^{-4}$
BCC-CSM1.1	2.6 [21.0]	2.5 [19.0]	1.3 [9.1]	0.7 [5.6]
CCM4	3.0 [16.0]	2.5 [19.0]	1.4 [14.0]	0.8 [8.6]
CNRM-CM5	3.5 [23.0]	2.9 [20.0]	1.9 [12.0]	0.7 [9.9]
GISS-E2-H	1.6 [8.8]	1.1 [7.2]	0.7 [5.6]	0.4 [3.4]
HadGEM-ES	3.8 [17.0]	3.1 [16.0]	2.0 [9.0]	0.4 [1.6]
IPSL-CM5a-LR	3.8 [24.0]	3.3 [27.0]	2.6 [20.0]	1.3 [13.0]
MPI-ESM-LR	3.6 [27.0]	2.8 [16.0]	1.5 [8.3]	0.8 [11.0]

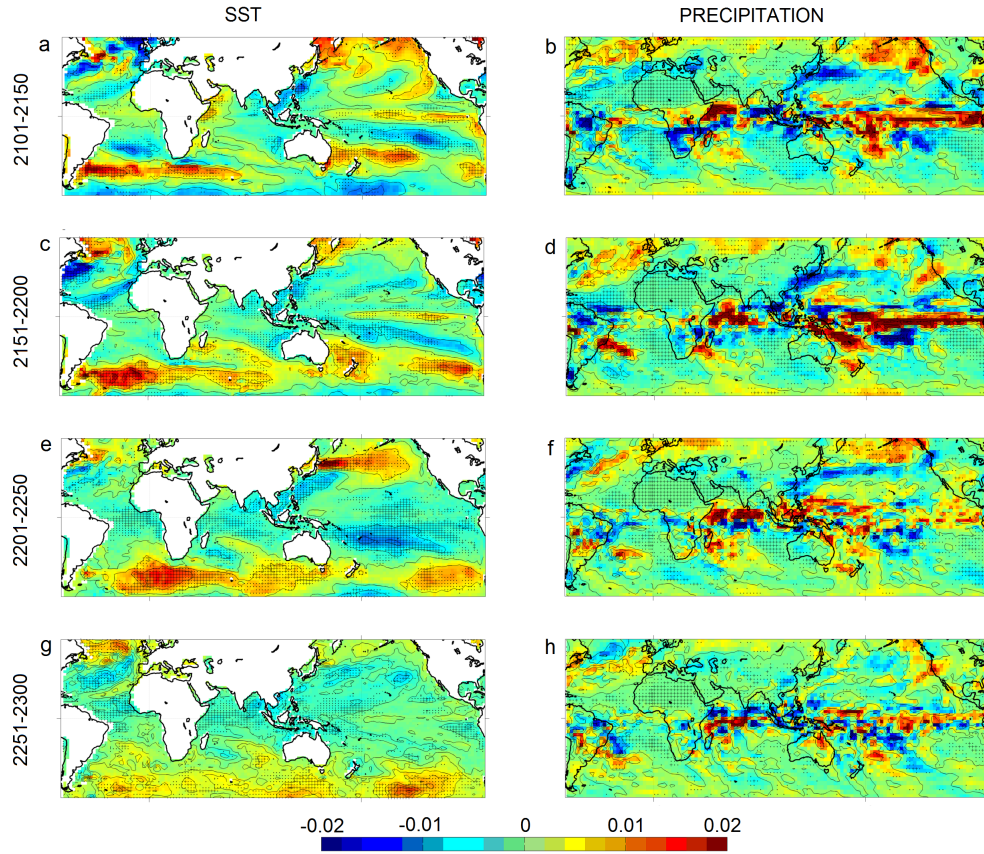


Figure 37: Trend anomaly maps for boreal winter in the 22nd and 23rd centuries. Anomalies are computed by removing the global mean trend calculated over the months of December to February and indicated in Table 4 from each grid cell. + and • indicate agreement in more than 90% and 70% of models in the sign of the trend anomaly slope. (a) Sea surface temperature (SST) averaged across models over 2101-2150. (b) Rainfall averaged across models over 2101-2150. (c) As in (a) but for 2151-2200. (d) As in (b) but for 2151-2200. (e) As in (a) but for 2201-2250. (f) As in (b) but for 2201-2250. (g) As in (a) but for 2251-2300. (h) As in (b) but for 2251-2300. The units are $^{\circ}\text{C}/\text{year}$ for SST and $(\text{mm}/\text{day})/\text{year}$ for precipitation

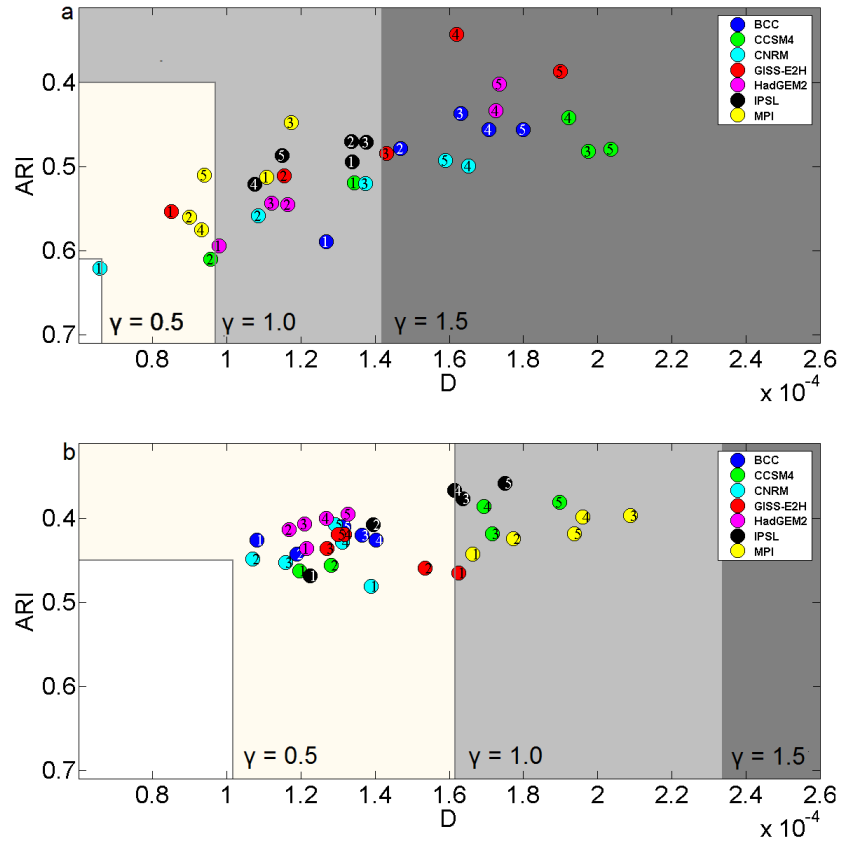


Figure 38: Metric D versus ARI for seven climate model networks from 2051 to 2300 over five consecutive 50-year periods, from 1 to 5. (a) Sea surface temperature. (b) Precipitation. All networks are referenced to the corresponding integration over the historical period. Three levels of noise-to-signal ratios γ are also indicated

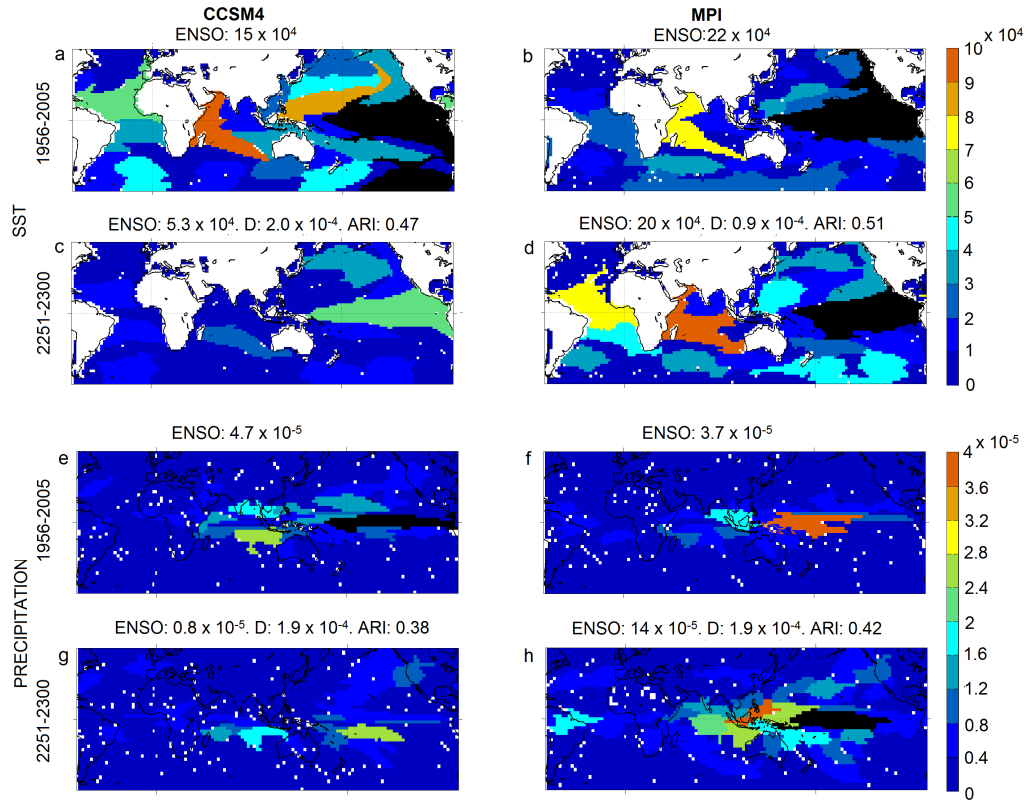


Figure 39: Sea surface temperature (a-d) and precipitation (e-h) strength maps for two models (left column CCSM4, right column MPI) in the historical period (1956-2005) and in the future (2251-2300). For each variable the first row corresponds to the historical experiments. For clarity, the strength of the ENSO-related area is saturated when exceeding the colorscale and indicated at the top of each panel. D and ARI metrics of the future projections from the corresponding historical member are also included

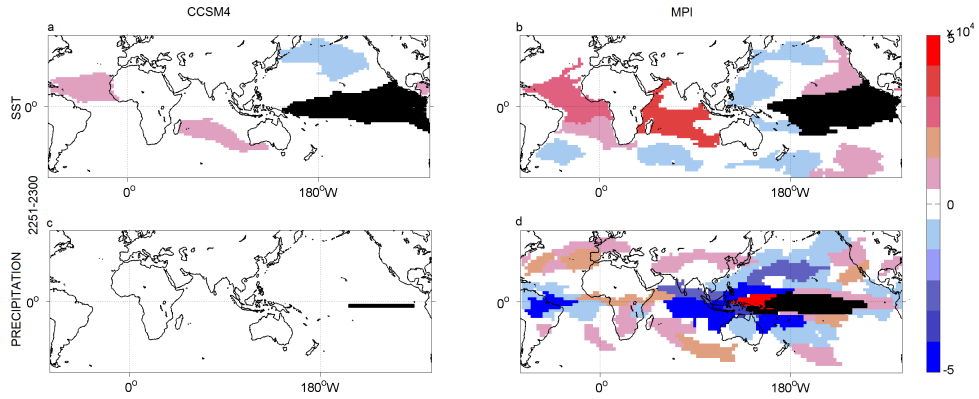


Figure 40: Link maps for sea surface temperature (a-b) and precipitation (c-d) from the ENSO-related area in black for two models for which the ENSO projected strength evolves in opposite ways. CCSM4 is shown on the left column and MPI on the right. Maps are calculated over the 2251-2300 period

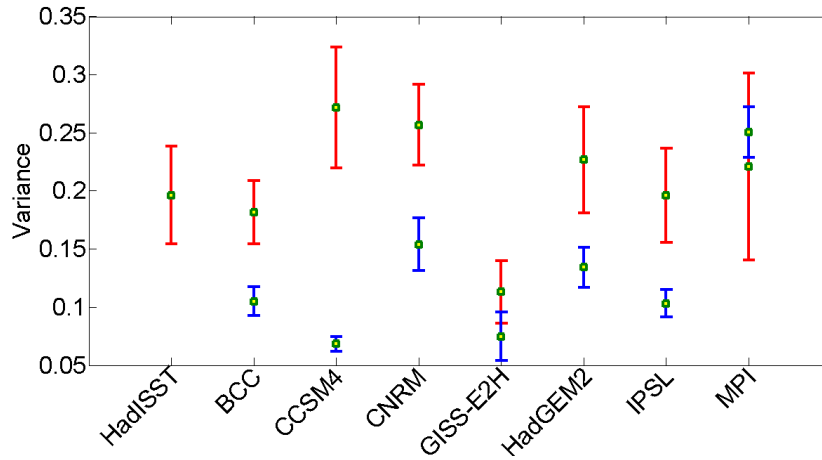


Figure 41: Variance of the cumulative anomalies of the ENSO area in DJF in the models and HadISST over 1956-2005 in red, and in the models over 2251-2250 in blue. For HadISST the time series is highly correlated (coefficient 0.94) with the Niño3.4 index defined as the average of SST anomalies from $5^{\circ}S$ to $5^{\circ}N$, and from 120° to $170^{\circ}W$. Error bars around the mean variance over 50 years are determined using a 20-year sliding window, and provide a measure of the decadal modulation of ENSO in the models over the periods considered.

4.4 Discussion

In this work we have established the stability of the SST and precipitation networks for twelve model ensembles in the CMIP5 catalog using a novel framework based on complex network analysis. This fast, scalable and robust method provides considerable advantages when comparing climate fields compared to more traditional approaches (e.g., predefined climate indices or EOFs). The areas identified reduce the dimensionality of the climate field and provide a compact and spatially embedded representation of the major modes of climate variability. Their interdependencies are quantified using weighted links, enabling us not only to detect their existence, but also to estimate their magnitude. With two metrics, the Adjusted Rand Index and a network distance metric, the output of climate models can be compactly validated against observations, intra-ensemble variability can be assessed, and networks obtained from model outputs under different forcing conditions can be contrasted. The applicability of the method is general and fits the objectives of any spatio-temporal data analysis, discovering unknown functional components and their inter-dependencies.

An important distinction between earlier network-based approaches and our method is that we construct networks that are complete and weighted graphs between homogeneous spatial areas. The clustering of grid cells into areas, the lack of edge pruning as well as the way in which we calculate the weights of the links between areas, makes the proposed network inference method more robust with respect to the underlying threshold compared to approaches that are based on (typically pruned) cell-level networks (as shown in Section 4.6).

The CMIP5 models have been validated against reanalyses over the second half of the 20th century, and compared for their projected responses under high GHG concentrations. We focused on global quantities, and analyzed fifty years time intervals; the dominant mode of variability at those time and spatial scales is ENSO, which induces the most severe global impacts in surface temperatures and precipitation, among other variables. Despite decades of research, ENSO sensitivity to changes in GHG concentrations remains undetermined in

the last generation of climate models [18, 39, 129].

The results of our analysis can be summarized as follows:

- Within the CMIP5 inventory, several models reproduce closely the observed SST network over the historical period (1956-2005), providing an accurate representation of major modes of climate variability and their links, despite biases in the climatologies. The spread in ARI and D between SST networks from members of the same ensemble is broadly consistent with the spread between different observational datasets or reanalyses. Precipitation networks, unsurprisingly, indicate that spatial likeness and strength are still challenging for modelers. However, the limited agreement between reanalysis products, and the evaluation of the noise-to-signal ratio suggests that the spatial and temporal intermittency of precipitation intrinsically limits the reproducibility of its topology. For rainfall, the intra-model spread in network metrics is generally very small; additionally, slope and patchiness of regional trends are decidedly underestimated by models. Together those outcomes suggest that CGCMs cannot yet capture the observed natural variability of rainfall. Models characterized by large D and small ARI in their SST fields, are also inaccurate in the representation of precipitation, but model performing the closest to the reanalysis in each of those fields differs.
- Changes in the network properties between the second half of the 20th and 21st centuries are generally modest and contained within the spread between different observational proxies in the historical period, despite substantial trends. This is especially true for the models that reproduce accurately the recent past. For those models uncertainties are greater in the projected trends than in the response of their modes of variability. Differences are slightly more probable in strength than in the spatial distributions of areas. Changes in distance D greater than 30% around the historical value signals the model tendency towards strengthening or weakening of major climate modes, and of ENSO, its variance, and its connectivity, even when

limited to one ensemble member. Eight of the twelve models analyzed display substantially weaker tropical areas and connections in one or more members, implying a decrease in ENSO strength and in potential predictability at seasonal and longer scales in the future [47]. The weakening of the links from the ENSO area for the majority of models analyzed is opposite to the conclusion presented in [29] and obtained characterizing El Nio activity using equatorial indices. Only two models, MPI and MIROC5, display a clear trend towards intensifying the strength of the ENSO area and its links.

- After 2100, models forced by the concentration pathway of the scenario with the highest greenhouse gas concentrations in CMIP5 reveal discernible changes in the strength of all major areas. Five out of seven follow an irreversible trajectory towards reducing dramatically the strength and size of the ENSO node, and towards weakening all ENSO links over the 23rd century. This behavior is mirrored in precipitation to a lesser extent. IPSL weakens as well, but partially recovers by 2300 in both SST and rainfall. MPI by the end of the integration has a virtually unaltered network in SST, while strength and links of the ENSO area increase substantially for precipitation. No obvious relation has been found between the trend patterns in the equatorial tropical Pacific, indicative of mean state changes, or the global warming/wetening trends, and the ENSO behavior in the networks [40], or between the response patterns of clouds and precipitation to a uniform warming in an aquaplanet configuration for three of the atmospheric components of the models analyzed [146], and their tropical rainfall response in a coupled set up. On the other hand, an increased tendency for eastward propagation of SST anomalies during positive ENSO events [129] in the 21st century, counterintuitively, may be symptomatic of an irreversible weakening of ENSO in the next century and of a loss of potential predictability in the atmosphere.

Considering the global impacts of tropical teleconnections and the changes in temperature and precipitation associated with El Niño and La Niña events (e.g., [9, 85]), we

conclude that the uncertainty in the projected connectivity of the climate system after 2100 in many regions and for models performing well under current conditions exceeds the uncertainty associated with the equilibrium temperature change.

For the question of robustness versus sensitivity of climate patterns under different forcing scenarios, the lack of consistency between models highlights, once more, the complexity associated with having multiple, nonlinear coupled processes. By adopting a perturbation-based approach and focusing on models with networks in the recent past that compare well to observations but diverge substantially in the future, it is possible to target more effectively efforts to understand the physical mechanisms and model parameterizations that cause such divergences.

4.5 Supplementary strength and link maps

Strength maps for boreal winter SST and precipitation for one member of each model considered are displayed below for the 1956-2005 historical period, the 2051-2100 RCP8.5 interval, and the ECP8.5 extension. Additionally, the corresponding link maps of the ENSO area are provided for both fields.

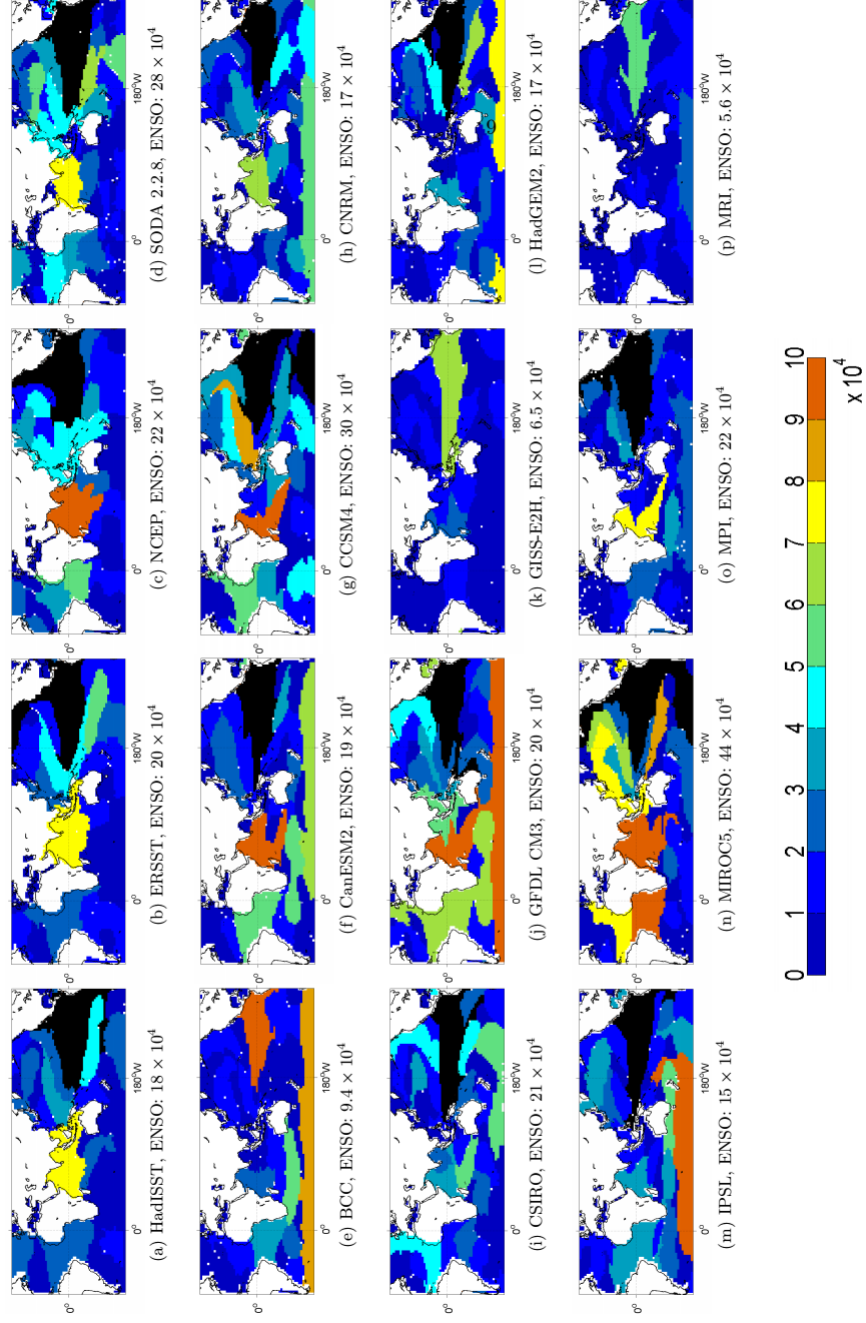


Figure 42: Maps of area strength of sea surface temperature networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded

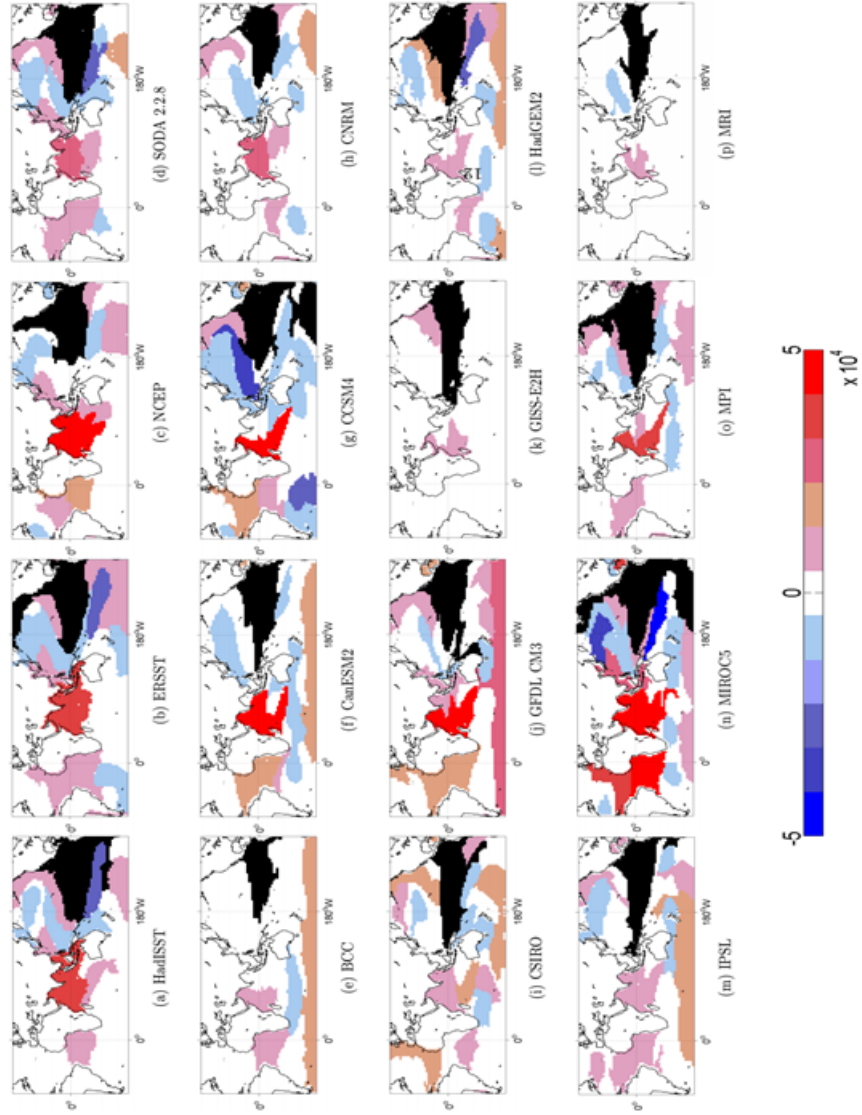


Figure 43: Link maps from the ENSO related area (in black) for sea surface temperature networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown

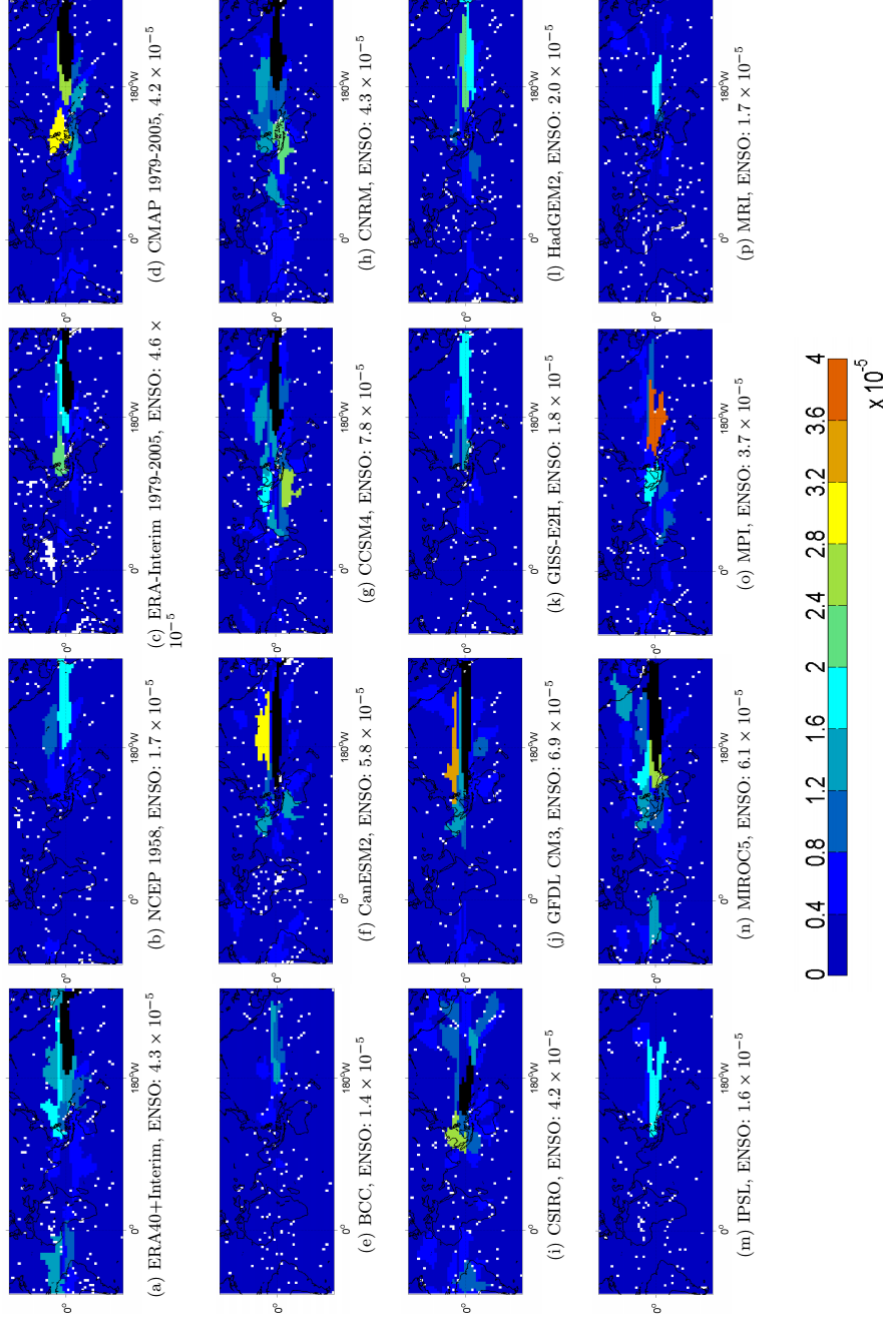


Figure 44: Maps of area strength for precipitation networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded

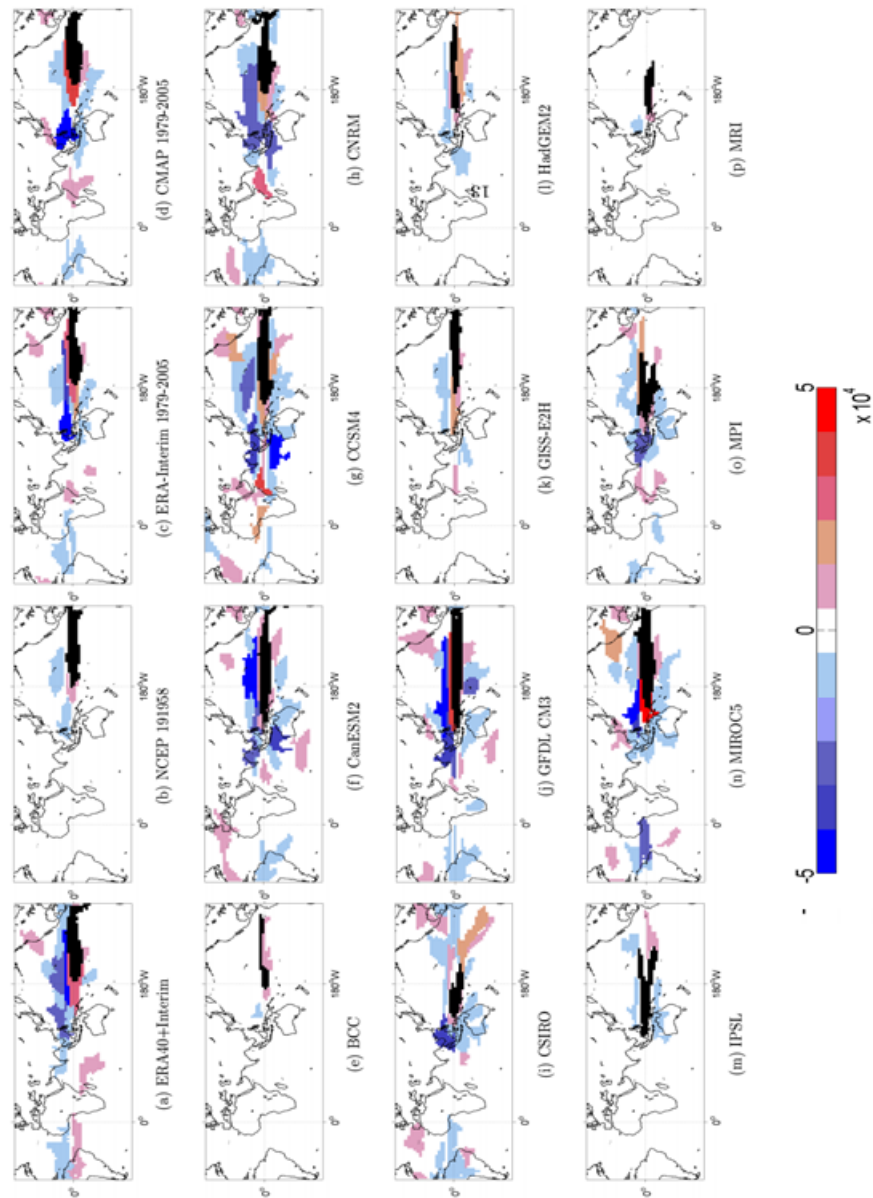


Figure 45: Link maps from the ENSO related area (in black) for precipitation networks in boreal winter (December to February) in the historical period 1956-2005 for models and reanalyses. Only one ensemble member per model is shown

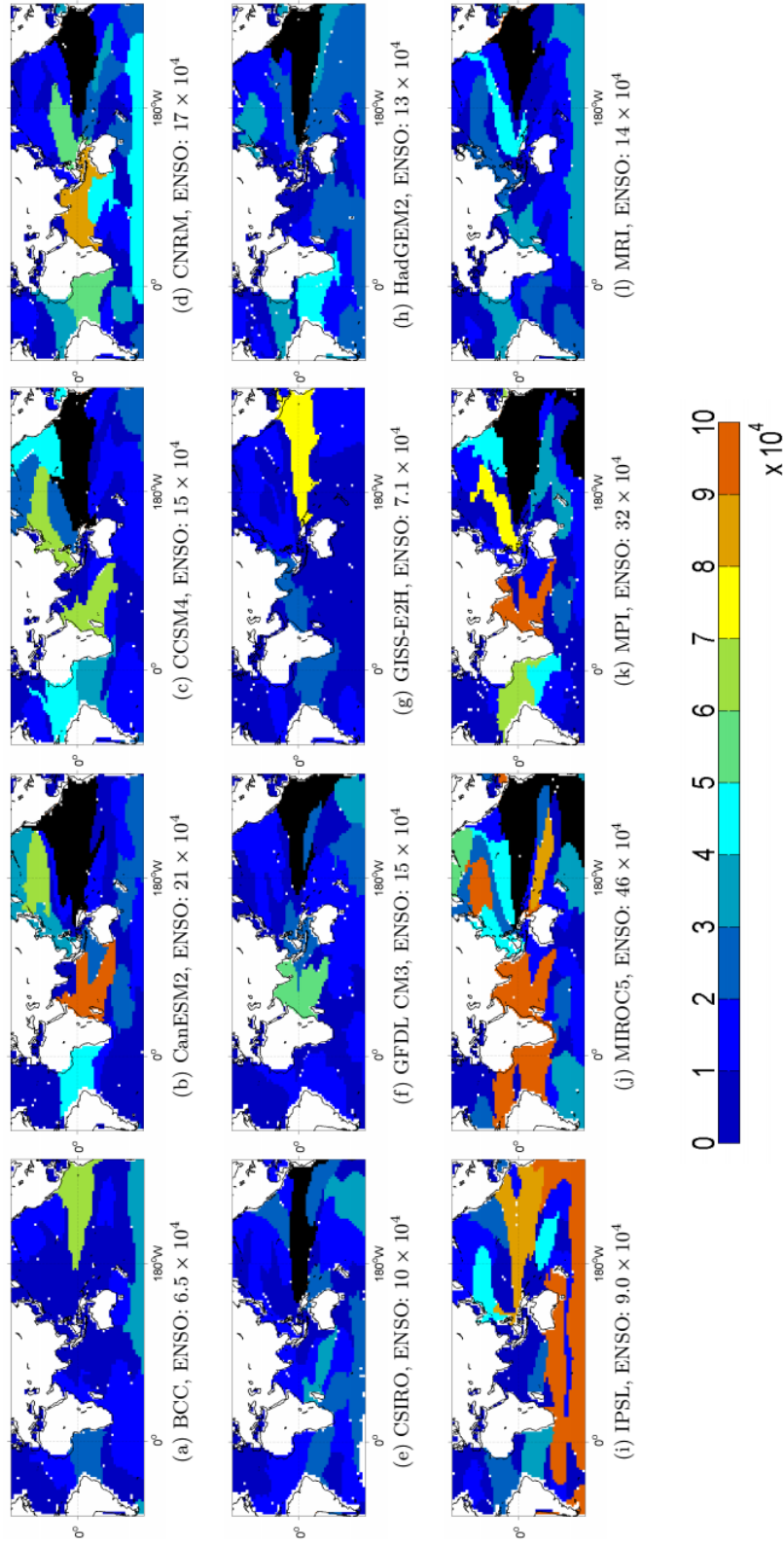


Figure 46: Maps of area strength for the sea surface temperature networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100). For each model, the ensemble member shown projects into the future the historical counterpart in Fig. 42. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded

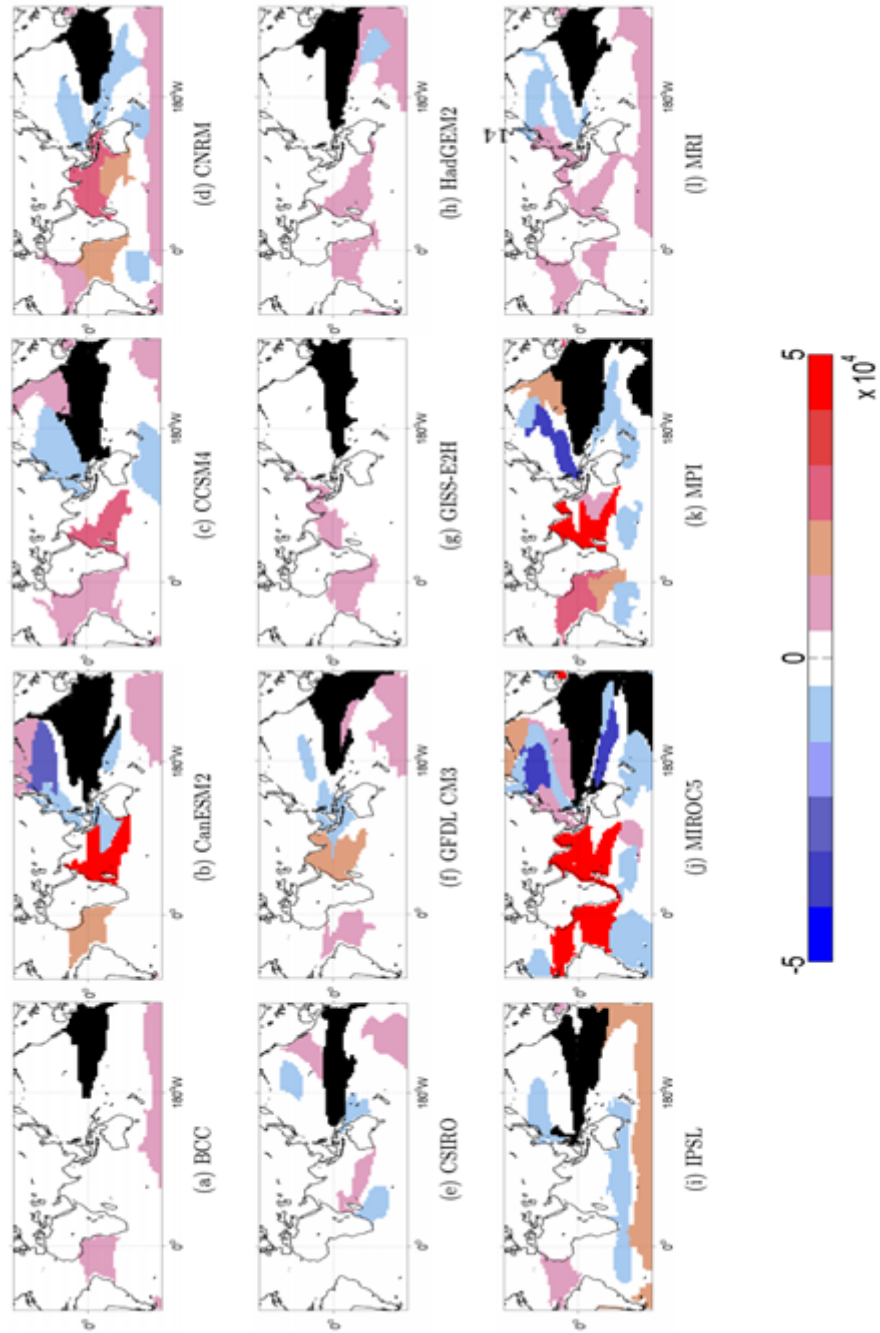


Figure 47: Link maps from the ENSO related area (in black) for sea surface temperature networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100) for the ensemble members in Fig. 46

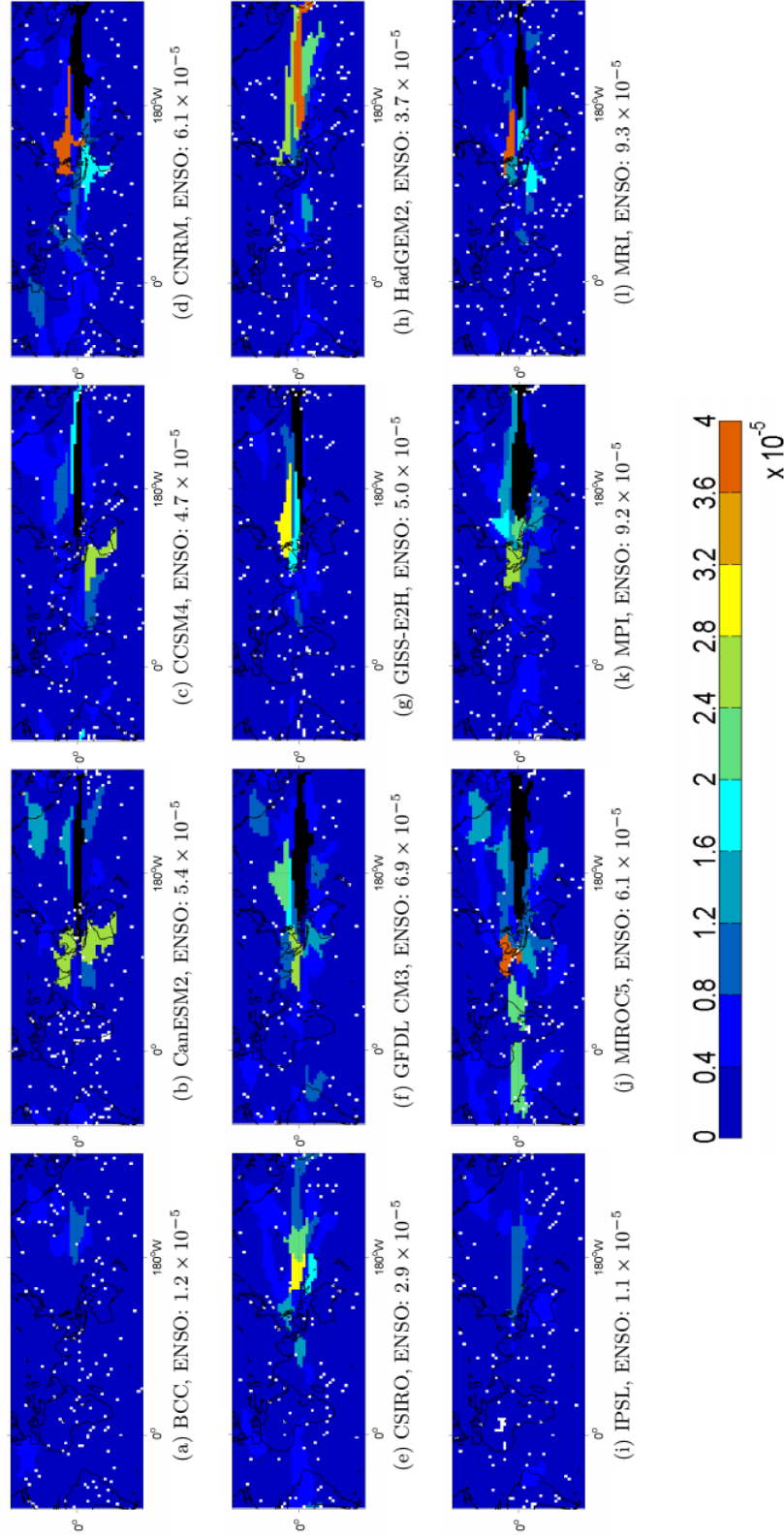


Figure 48: Maps of area strength for the precipitation networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100). For each model, the ensemble member shown is the projection into the future of the historical counterpart in Fig. 44. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded

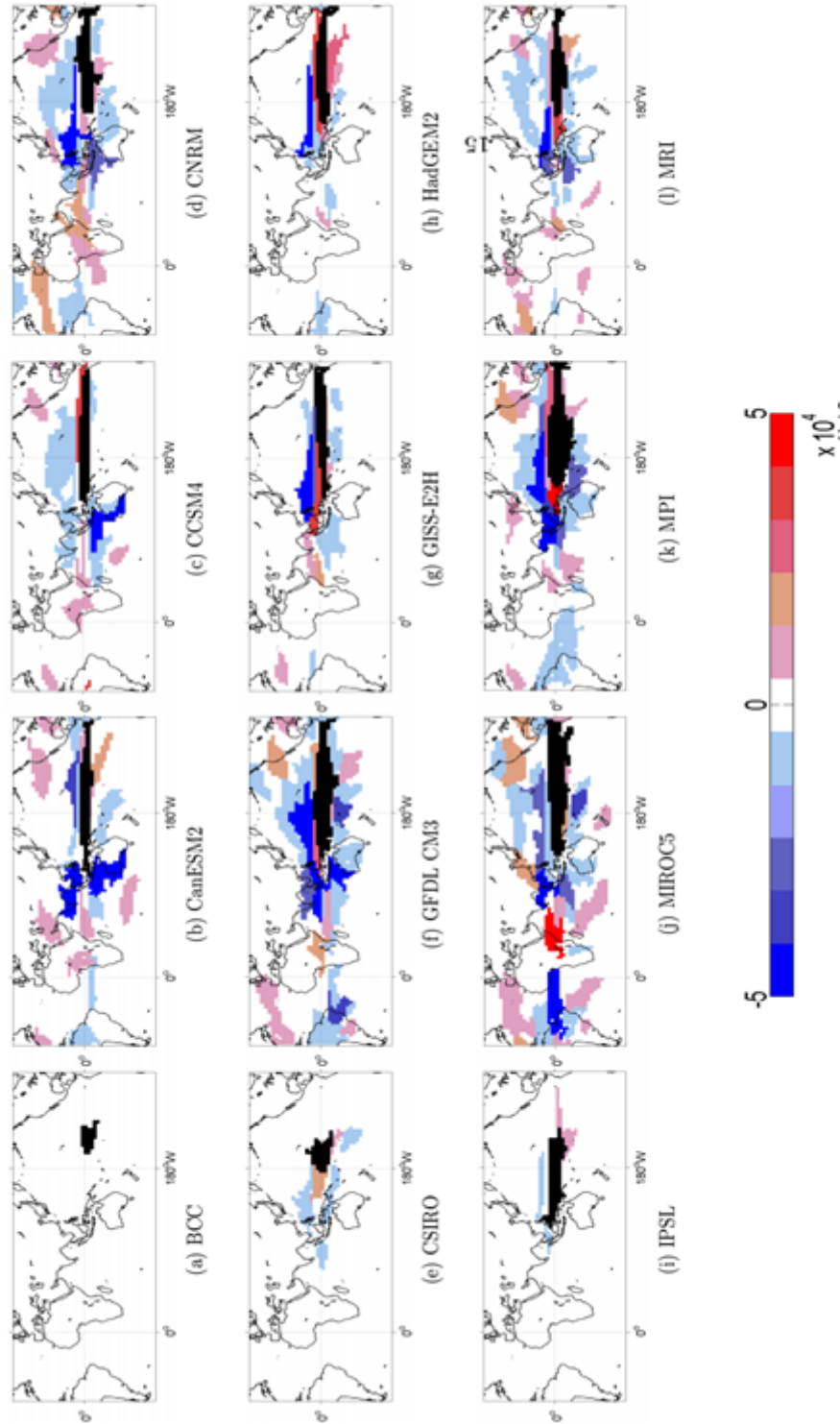


Figure 49: Link maps from the ENSO related area (pictured in black) for the precipitation networks in boreal winter (December to February) in the RCP8.5 projections (period 2051-2100) for the ensemble members in Fig. 48

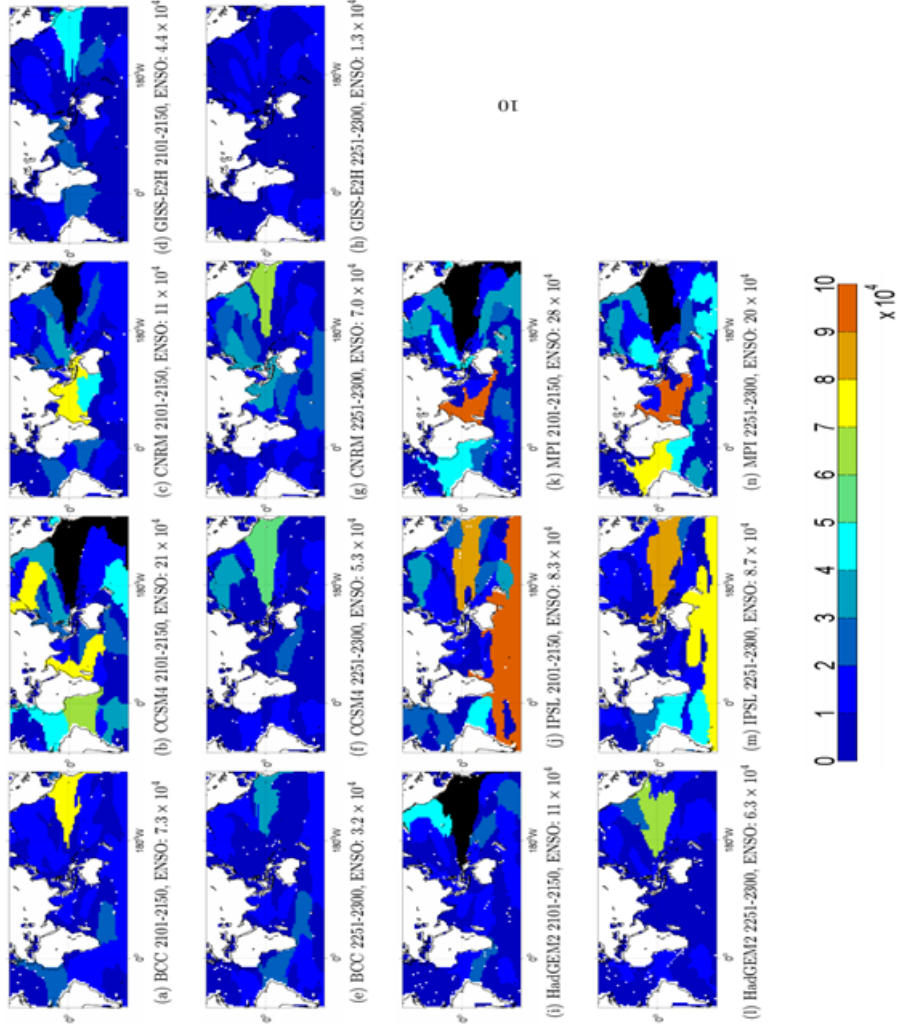


Figure 50: Maps of area strength for the sea surface temperature networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300). For each model, the ensemble member shown projects into the future the historical counterpart in Fig. 42. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded

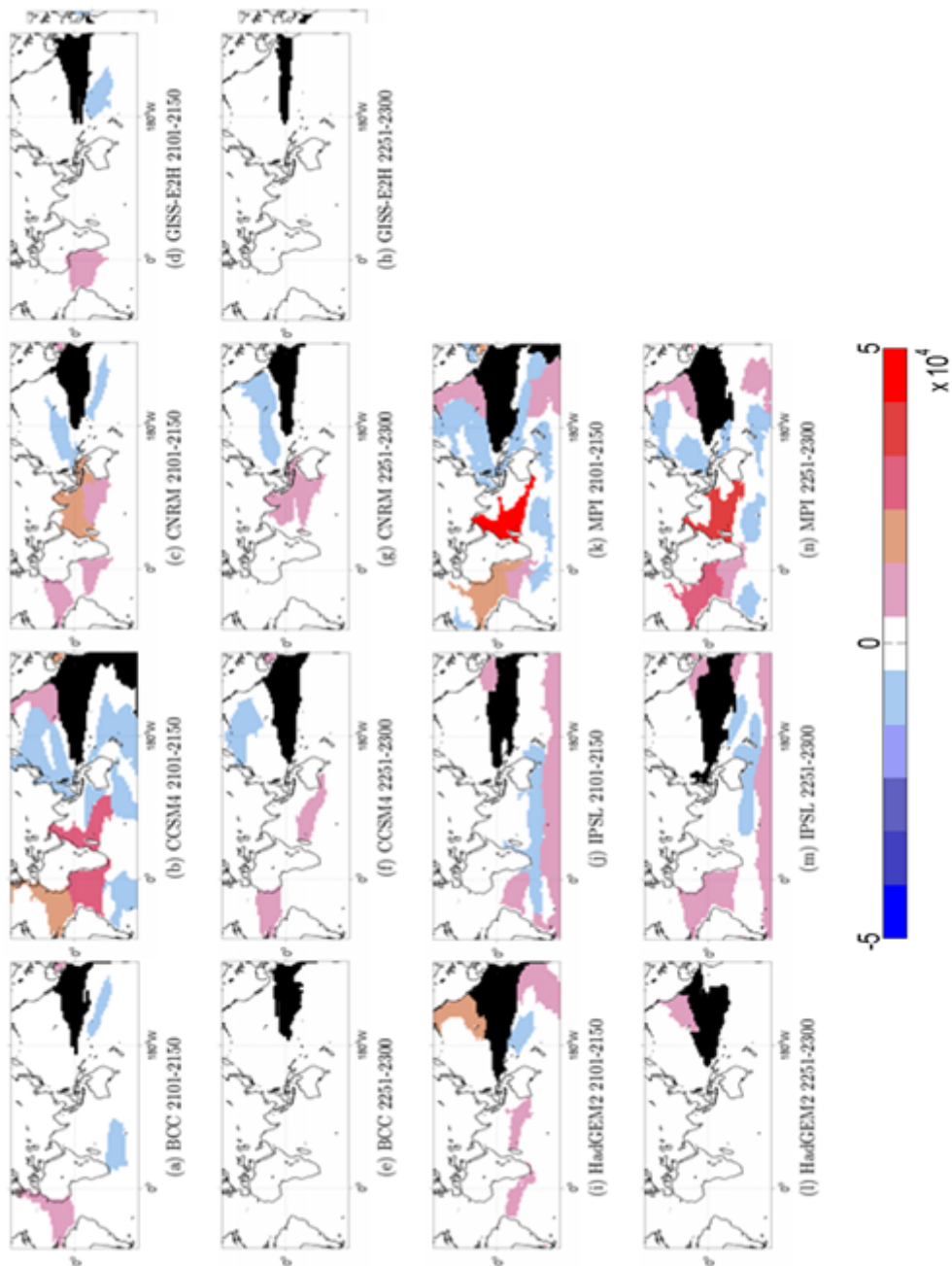


Figure 51: Link maps from the ENSO related area (in black) for sea surface temperature networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300) for the ensemble members in Fig. 50

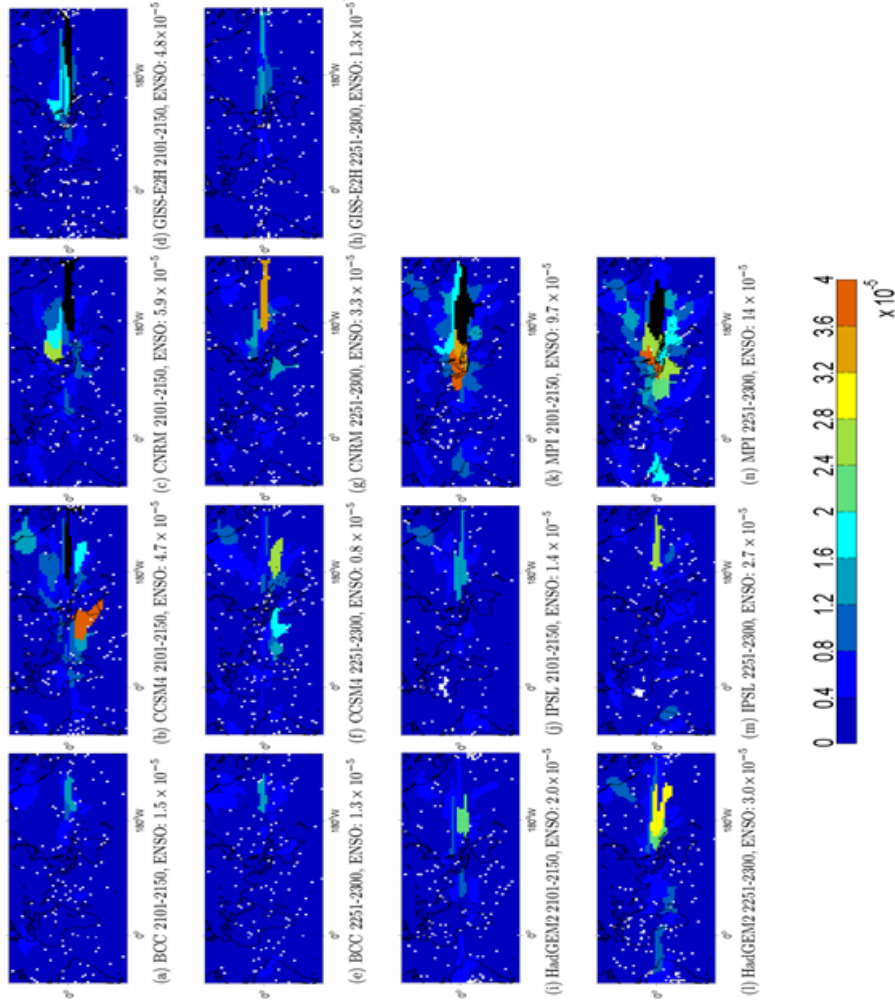


Figure 52: Maps of area strength for the precipitation networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300). For each model, the ensemble member shown projects into the future the historical counterpart in Fig. 44. The strength of the area corresponding to ENSO is indicated in the panel captions and saturated in black if colorbar limits are exceeded

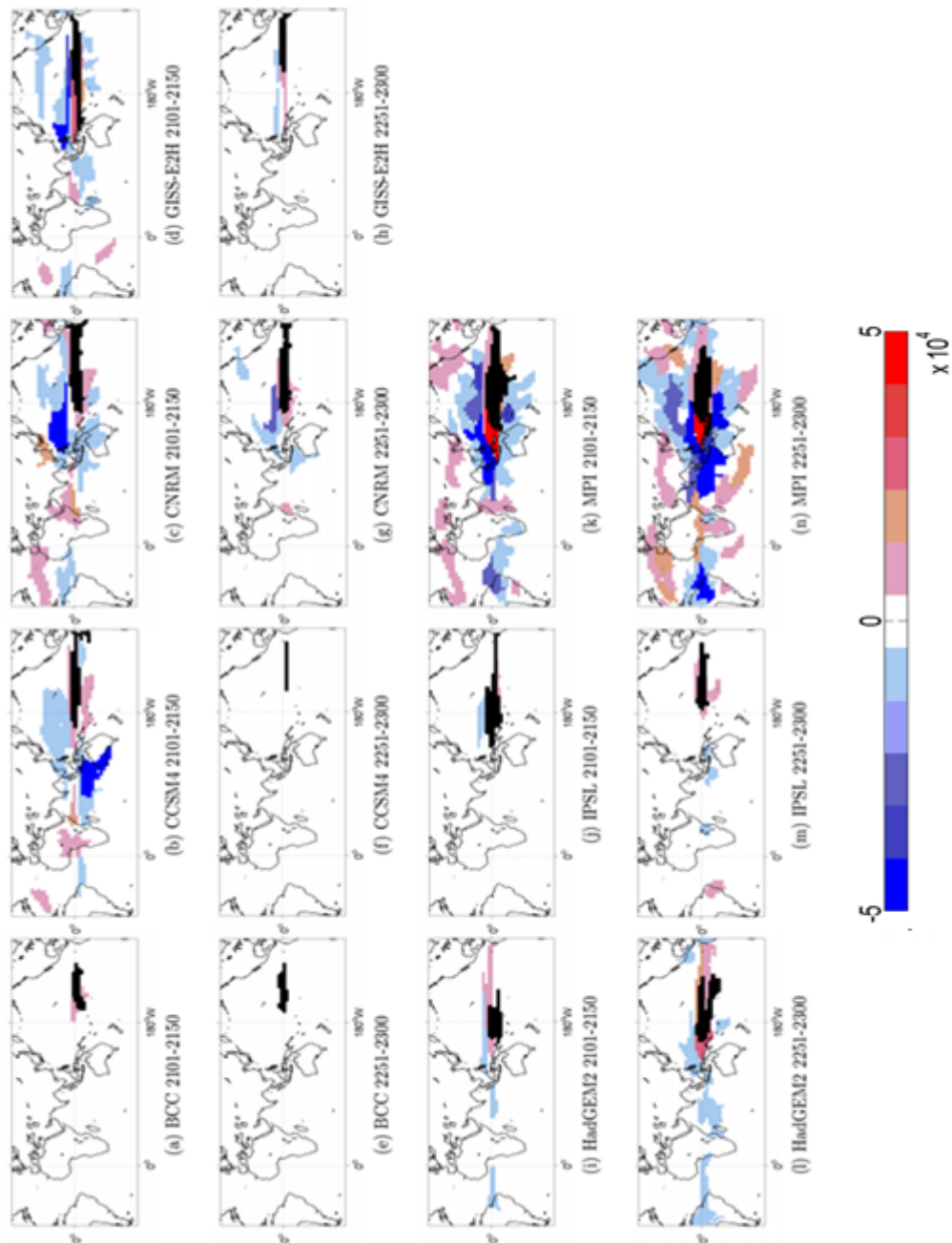


Figure 53: Link maps from the ENSO projections (in black) for the precipitation networks in boreal winter (December to February) in the ECP8.5 projections (periods 2101-2150 and 2251-2300) for the ensemble members in Fig. 52

4.6 Advantages of using a complete weighted cell-level network

The network methodology adopted in this work differs from several others proposed in the literature [143, 153] most importantly because no pruning is done to remove edges. Here both the cell-level network and the area-level networks are modeled as complete weighted graphs. Instead, the majority of earlier climate network inference methods construct unweighted graphs in which whenever the cross-correlation between two cells is less than a threshold, the corresponding cells are not considered connected.

Our approach offers two substantial advantages. First, by modeling the climate network as a weighted graph we can leverage information about the actual magnitude of the cell-level correlations. The information captured by these weights can give us insights about the strength of specific teleconnections between different nodes of the network (e.g. between ENSO and areas forming the horseshoe pattern). Secondly, the proposed method is more robust compared to methods that perform pruning. Robustness is an important property, especially for the objectives of this chapter, since we compare different climate models and the properties of the climate system they simulate over time.

In this section, we substantiate those points by showing that link pruning makes the network inference process less robust, based on two comparisons. Our proposed inference method relies on a single parameter, the level of significance α . The parameter τ , which is used in the area detection algorithm, is calculated based on α , as described in Section 3.8. Let r_α denote the minimum significant correlation for a given level of significance α . In the first comparison the input to the area identification algorithm is an unweighted network. All pair-wise grid cell correlations that are non-significant for the given level α are set to 0. Correlations larger than r_α are set to 1 and correlations lower than $-r_\alpha$ are set to -1. We refer to the corresponding cell-level network as the unweighted pruned network (such networks have been studied in [61, 145]). The second comparison is a more relaxed version of the first; we simply remove all pair-wise cell correlations that are non-significant for a given level of significance α but maintain the actual magnitude of the significant

correlations. We refer to this type of cell-level network as the weighted pruned network (such networks have been studied in [143]). In both cases, the τ threshold is computed as in the weighted complete network. Consequently, the area identification algorithm and the threshold τ are the same for all three networks; the only difference is the input to the area identification algorithm (i.e. the cell-level network).

In all comparisons our "reference network" is constructed using the HadISST 1956-2005 (DJF) anomaly time series for $\alpha = 1 \times 10^{-3}$. The identified areas are presented in Fig. 54. When the input is the unweighted pruned network the resulting areas cannot be easily interpreted in a climate context. For example, the area corresponding to the Indian Ocean extends to the North Pacific Ocean while the ENSO related area includes ample extratropical regions in both hemispheres. The areas identified using a weighted pruned network are closer to those identified using a complete weighted graph but, as we shall prove next, the former is less robust. We cannot be certain that a certain set of areas is the "right set", since no ground truth exists for such an evaluation, but any network methodology used for model intercomparison should, at least, be robust to its input parameter (α in our case) and should be insensitive (or have only small sensitivity) to changes in α . To evaluate the robustness of the various methodologies we vary α around its standard value and quantify the network changes in terms of the ARI metric (Fig. 55). The network inference process is more robust when the cell-level network is modeled as a complete and weighted graph. If we prune some edges (and keep the weight of the remaining links) the robustness of the method decreases resulting in lower ARI values. Finally, the least robust option is to model the cell-level network as an unweighted pruned network.

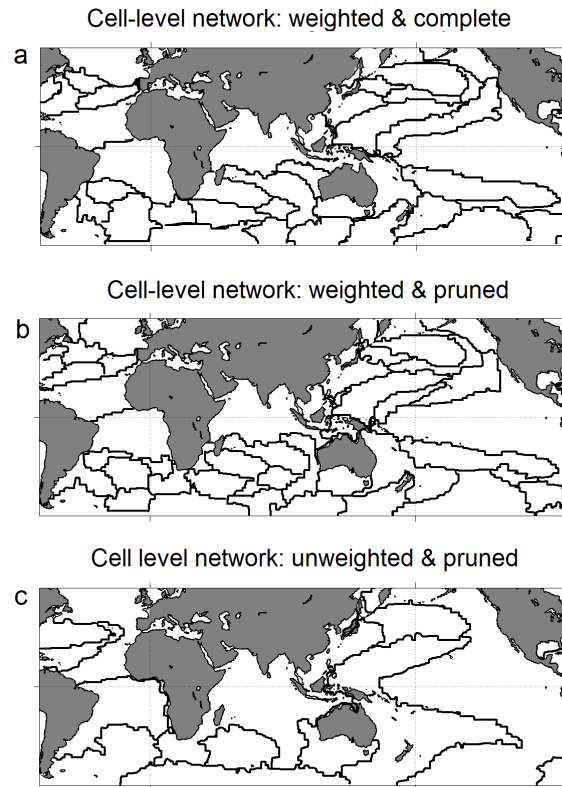


Figure 54: Areas identified using three different cell-level networks. α was set to 1×10^{-3} . Data set: HadiSST 1956-2005

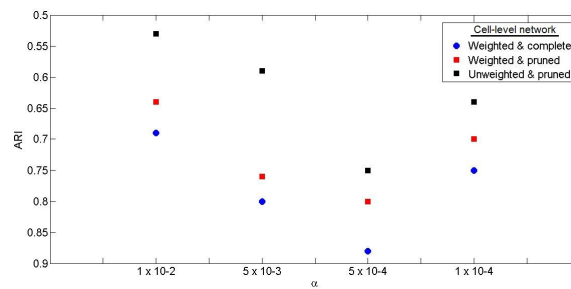


Figure 55: ARI between a reference network constructed using $\alpha = 1 \times 10^{-3}$ and networks constructed using different α values

Chapter V

δ -MAPS: FROM SPATIO-TEMPORAL DATA TO A WEIGHTED AND LAGGED NETWORK BETWEEN FUNCTIONAL DOMAINS

5.1 Introduction

Spatio-temporal data become increasingly prevalent and important for both science (e.g., climate, systems neuroscience, seismology) and enterprises (e.g., the analysis of geotagged social media activity). The spatial scale of the available data is often determined by an arbitrary grid, which is typically larger than the true dimensionality of the underlying system. One major task is to identify the distinct semi-autonomous components of this system and to infer their (potentially lagged and weighted) interconnections from the available spatio-temporal data. Traditional dimensionality reduction methods, such as PCA, ICA or clustering, have been successfully used for many years but they have known limitations when the objective is to infer the functional network between all spatial components of the system.

We propose δ -MAPS, an inference method that first identifies these spatial components, referred to as “domains”, and then the connections between them (§5.3). Informally, a *functional domain* (or simply *domain*) is a spatially contiguous region that somehow participates in the same dynamic effect or function. The exact mechanism that creates this effect or function varies across application domains; however, the key idea is that *the functional relation between the grid cells of domain results in highly correlated temporal activity*. If we accept this premise, it follows that we should be able to identify the “epi-center” or *core of a domain* as a point (or subregion) at which the local homogeneity is maximum across the entire domain. Instead of searching for the discrete boundary of a domain, which may not exist in reality, we compute a domain as the *maximum possible set*

of spatially contiguous cells that include the detected core, and that satisfy a homogeneity constraint, expressed in terms of the average pairwise cross-correlation across all cells in the domain. Domains may be spatially overlapping. Also, some cells may not belong to any domain.

After we identify all domains, δ -MAPS infers a functional network between them. Different domains may have correlated activity, potentially at a lag, because of direct or indirect interactions. The proposed edge inference method examines the statistical significance of each lagged cross-correlation between two domains, applies a multiple-testing process to control the rate of false positives, infers a range of potential lag values for each edge, and assigns a weight to each edge based on the covariance of the corresponding two domains.

δ -MAPS is related to clustering, parcellation (or regionalization), network community detection, multivariate statistical methods for dimensionality reduction such as PCA and ICA, as well as functional network and lag inference methods. However, as we discuss in §5.2 and show with synthetic data experiments in §5.4, δ -MAPS is also significantly different than all these methods. δ -MAPS does not require the number of domains as an input parameter, the resulting domains are spatially contiguous and potentially overlapping, and the inferred connections between domains can be lagged and positively or negatively weighted. Further, the distinction between grid cells that are correlated within the same domain and grid cells that are correlated across two distinct domains allows δ -MAPS to separate between local diffusion (or dispersion) phenomena and remote interactions that may be due to underlying structural connections (e.g., a white-matter fiber between two brain regions).

We illustrate the application of δ -MAPS on data from two domains: climate science (§5.5) and neuroscience (§5.6). First, the sea-surface temperature (SST) climate network identifies some well-known climate “tele-connections” (such as the lagged connection between the El Niño Southern Oscillation and the Indian ocean). Second, the analysis of resting-state fMRI cortical data confirms the presence of three well-known functional brain

“networks” (default-mode, occipital, and motor/somatosensory), and shows that the cortical network includes a *backbone* of relatively few regions that are densely interconnected.

5.2 Related Work

A common approach to reduce the dimensionality of spatio-temporal data is to apply PCA (standard or rotated) or ICA techniques. For instance, in climate science, PCA (also known as Empirical Orthogonal Function (EOF) analysis) has been used to identify teleconnections between distinct climate regions [167]. The orthogonality between PCA components complicates the interpretation of the results making it difficult to identify the distinct underlying modes of variability and to separate their effects, as clearly discussed in [50]. ICA analysis is more common in the neuroscience literature, aiming to identify independent rather than orthogonal components [88]. However, ICA does not provide a relative significance for each component, and the number of independent components should be chosen based on some additional information about the underlying system.

Another broad family of spatio-temporal dimensionality reduction methods is based on clustering [22, 60, 139, 177]. These algorithms can be grouped into region-growing methods (e.g., [23, 104]), spectral (e.g., the NCUT method often applied in fMRI analysis [44, 160] – but also see a discussion of their limitations [14]), hierarchical (e.g., [24, 150]), and probabilistic (e.g., [14, 83]). These groups of algorithms are quite different but they share some common characteristics: the resulting clusters may not be spatially contiguous, they are typically non-overlapping, every grid cell needs to belong to a cluster (potentially excluding only outliers), and the number of clusters is often required as an input parameter. In particular, the lack of spatial contiguity makes it hard to distinguish between correlations due to spatial diffusion (or dispersion) phenomena from correlations that are due to remote (structural) interactions between distinct effects.

An approach of increasing popularity is to first construct a correlation-based network

between individual grid cells, after pruning cross-correlations that are not statistically significant – see [100]. Then, some of these methods analyze the (binary or weighted) cell-level network directly based on various centrality metrics, k-core decomposition, spectral analysis, etc. (e.g., [52, 161]) or they first apply a community detection algorithm (potentially able to detect overlapping communities, e.g., [4, 103, 116]) on the cell-level network and then analyze the resulting communities in terms of size, density, location, overlap, etc. (e.g., [108, 118, 142, 143]). A community however may group together two regions that are, first, not spatially contiguous, and second, different in terms of how they are connected to other regions; an instance of this issue is illustrated in Fig. 58-C in the context of climate data analysis.

5.3 δ -MAPS

The input data is generated from a *spatial field* $\mathbf{X}(t)$ sampled on an arbitrary *grid* G . This grid can be modeled as a planar graph $G(V, E)$, where each vertex in V is a grid cell and each edge in E represents the spatial adjacency between two neighboring cells. A set of cells $A \subseteq V$ is *spatially contiguous*, denoted by $I_G(A)=1$, if it forms a connected component in G .

The K -*neighborhood* of a cell i , denoted by $\Gamma_K(i)$, includes i and the set of K nearest neighbors to i according to an appropriate spatial distance metric (e.g., geodesic distance for climate data, Euclidean distance for fMRI data). The K -neighborhood of a cell is always spatially contiguous.

Each grid cell i is associated with a time series $x_i(t)$ of length T ($t \in \{1, \dots, T\}$). We assume that $x_i(t)$ is sampled from a stationary signal and denote by $\tilde{\mu}_i$ and $\tilde{\sigma}_i^2$ its sample mean and variance, respectively. The similarity between the activity of two cells i and j is measured with Pearson's cross-correlation at zero-lag,

$$r_{i,j} = \frac{\sum_{t=1}^T (x_i(t) - \tilde{\mu}_i)(x_j(t) - \tilde{\mu}_j)}{T \tilde{\sigma}_i \tilde{\sigma}_j}. \quad (10)$$

Other similarity metrics could be used instead.

The *local homogeneity at cell i* is defined as the average pairwise cross-correlation between the $K + 1$ cells in $\Gamma_K(i)$,

$$\hat{r}_K(i) = \frac{\sum_{m \neq n \in \Gamma_K(i)} r_{m,n}}{K(K+1)}. \quad (11)$$

Similarly, we define the *homogeneity of a set of cells A* as the average pairwise cross-correlation between all distinct cells in A ,

$$\hat{r}(A) = \frac{\sum_{m \neq n \in A} r_{m,n}}{|A|(|A| - 1)}. \quad (12)$$

5.3.1 Functional domains

Intuitively, a *domain A* is a spatially contiguous set of cells that somehow participate in the same dynamic effect or function. The exact mechanism that creates this effect or function varies across application domains; however, the key premise is that *the functional relation between the cells of domain A results in highly correlated temporal activity (at zero-lag), and thus high values of the homogeneity metric $\hat{r}(A)$* . A given *homogeneity threshold δ* examines if the homogeneity of A is sufficiently high, i.e., a domain A must have $\hat{r}(A) > \delta$. (the selection of δ is discussed later in this section).

If we accept this premise, it follows that we should be able to identify the “epicenter” or *core of a domain A* as a cell $i \in A$ at which the local homogeneity $\hat{r}_K(i)$ is maximum across all cells in A (and certainly larger than δ). In general, the core of a domain may not be a unique cell.

More formally now, suppose that we know that cell c is in the core of a domain. The *domain A rooted at c* has to satisfy the following three properties: it should include cell c , be spatially contiguous, and have higher homogeneity than δ :

$$c \in A, \quad I_G(A) = 1, \quad \hat{r}(A) > \delta. \quad (13)$$

A domain may not have sharp spatial boundaries; instead, it may gradually “fade” into other domains or regions dominated by noise. So, instead of searching for the discrete

boundary of a domain, it is more reasonable to compute a domain as the *largest possible set of cells* that satisfies the previous three constraints.

Domain identification problem: Given the field $\mathbf{X}(t)$ on the spatial grid G , a core cell c , and the threshold δ , the domain $A(c)$ is a maximum-sized set of cells that satisfies the three constraints of (13). In Section 5.8 we prove that the decision version of this problem is NP-Hard.

A given spatial field $\mathbf{X}(t)$ may include several domains. The number of identified domains, denoted by N , depends on the threshold δ . Domains may be spatially overlapping; this is the case when the cells of a region are significantly correlated with two or more distinct domain cores. Also, some cells of the grid may not belong to any domain, meaning that their signal can be thought of as mostly noise (at least for the given value of δ). Decreasing δ will typically result in a larger number of detected domain cores. Further, as δ decreases, the spatial extent of each domain will typically increase, resulting in larger overlaps between nearby domains.

δ can simply be a user-specified parameter for the minimum required average cross-correlation within a domain. Another way is to calculate δ based on a statistical test for the significance of the observed zero-lag cross-correlations. A summary of this method is given next (described in more detail in Section 5.9). We start with a random sample of pairs of grid cells. We then apply the statistical test described in §5.3.2 (see Equations 15 and 16) to examine if the zero-lag cross-correlation between each of these pairs passes a given significance level α (set to 10^{-2} unless specified otherwise). δ is then set to the average of the statistically significant cross-correlations in that sample. The rationale is that the average pairwise cross-correlation among cells that belong to the same domain should be higher than a sample average of statistically significant cross-correlations between cells that can be anywhere on the grid.

5.3.1.1 Algorithm for domain identification

Given the NP-Hardness of the previous problem, we propose a greedy algorithm that runs in two phases. In the first phase, we identify a set of cells, referred to as *seeds*; each seed is a candidate core for a domain. In the second phase, each seed is initially considered as a distinct domain. Then, an iterative and greedy algorithm attempts to identify the largest possible domains that satisfy the three constraints of (13) through a sequence of *expansion* and *merging* operations. The two phases are described next, while the complete pseudocode is presented in Section 5.10. The source code (including supporting documentation) will be available on GitHub.

Seed selection Recall that the core of a domain is a cell of maximum local homogeneity across all cells of that domain. So, one way to detect *potential* core cells, while the domains are still unknown, is to identify points at which the homogeneity field $\hat{r}_K(i)$ is locally maximum. Specifically, cell i is a seed if $\hat{r}_K(i) > \delta$ and $\hat{r}_K(i) \geq \hat{r}_K(j) \forall j \in \Gamma_K(i)$. Let S be the set of all identified seeds.

In general, a single domain may produce more than one seed because the local homogeneity field can be noisy and so it may include multiple local maxima, greater than δ . Further, additional seeds can appear in regions where domains overlap. Consequently, it is necessary to include a merging operation in which two or more seeds are eventually merged into the same domain.

Note that as K decreases, the local homogeneity field becomes more noisy and so we may detect more seeds in the same domain. On the other hand, larger values of the neighborhood size K can oversmooth the homogeneity field, removing seeds and potentially hiding entire domains. The latter is more likely if the spatial extent of a domain is smaller than $K+1$ cells. This observation implies that the spatial resolution of the given grid sets a lower bound on the size of the functional domains that can be detected.

Domain-merging operation Two candidate domains A and B can be merged if they are spatially contiguous and if the homogeneity of their union is sufficiently high, i.e., $\hat{r}(A \cup B) > \delta$. Whenever there is more than one pair of domains that can be merged, we greedily choose the pair with the maximum union homogeneity; this greedy choice makes the merged domain more likely to expand further.

The merging operation is performed initially on the set of seeds S . It is also performed after each domain-expansion operation, whenever it is possible to do so.

Domain-expansion operation A domain A is expanded by considering all cells that are adjacent to A , and selecting the cell i that maximizes $\hat{r}(A \cup \{i\})$; again, this greedy choice makes the expanded domain more likely to expand further.

The expansion operation is repeated in rounds. At the start of each round, domains are sorted in decreasing order of homogeneity. Then, each domain is expanded by one cell at a time, as previously described, in that order. After every expansion operation, we check whether one or more merging operations are possible. A round is complete when we have attempted to expand each domain once.

A domain can no longer expand if that would violate the homogeneity constraint δ or if there are no other adjacent cells that can be added into the domain. The domain identification algorithm terminates when no further expansion or merging operations are possible.

5.3.2 The domain network

Given the N identified domains $V_\delta = \{A_1, \dots, A_N\}$, the next step is to construct a network $G_\delta(V_\delta, E_\delta)$ between domains. Different domains may have correlated activity because of direct or indirect interactions. We refer to G_δ as a *functional network* to emphasize that the edges between domains are based on functional activity and correlations instead of structural or physical connections (“structural network”) or causal interactions (“effective network”).

We associate a *domain-level signal* $X_A(t)$ with each domain A . The definition of this signal depends on the specific application field. For instance, when we analyze climate anomaly time series, the domain-level signal is defined as the *cumulative anomaly* across all cells of that domain, where the contribution of each signal is weighted by the relative size of that cell (it depends on the cell's latitude). For fMRI data, the domain-level signal is defined as the *average BOLD signal* across the cells of that domain.

Two different domains may be located at some distance, and so they may be correlated at a non-zero lag τ . For this reason, we examine if there is a significant cross-correlation between different domains over a range of lags ($-\tau_{max} \leq \tau \leq \tau_{max}$). The sample cross-correlation between domains A and B at a lag τ can be estimated as:

$$r_{A,B}(\tau) = \frac{\sum_{t=1}^{T-\tau} (X_A(t) - \tilde{\mu}_A)(X_B(t + \tau) - \tilde{\mu}_B)}{T \tilde{\sigma}_A \tilde{\sigma}_B}, \quad (14)$$

where $\tilde{\mu}_A$ and $\tilde{\sigma}_A$ denote sample mean and standard deviation estimates, respectively. The selection of τ_{max} should be large enough to include the typical signal propagation delays in the underlying system but at the same time it should be much lower than T . The $2\tau_{max} + 1$ cross-correlations for a pair of domains can be represented with a *correlogram*; an example based on climate sea-surface temperature data (see §5.5) is shown in Fig. 56.

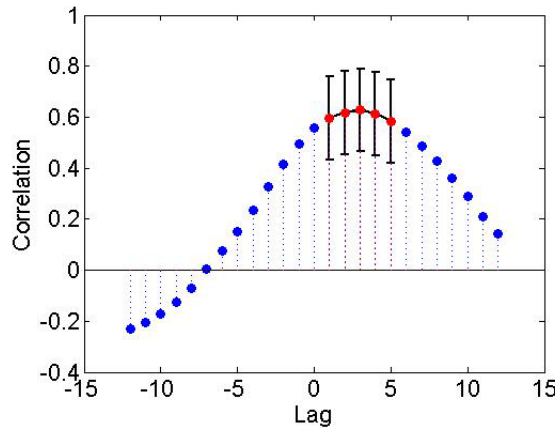


Figure 56: Correlogram between two climate time series for a lag range of ± 12 months. We show the significant correlations for a false discovery rate $q = 10^{-3}$ with red. The error bars correspond to \pm one standard deviation, as estimated by Eq. (15).

The next step is to examine the statistical significance of the measured cross-correlation between two domains A and B . Two uncorrelated signals can still produce a considerable sample cross-correlation if they have a strong auto-correlation structure. This is captured by Bartlett's formula [26], which is an estimator for the variance of $r_{A,B}(\tau)$ (for a fixed value of τ). Under the null-hypothesis that the domain-level signals of A and B are uncorrelated,

$$\text{Var}[r_{A,B}(\tau)] = \frac{1}{T - \tau} \sum_{\tau_k = -T}^T r_{A,A}(\tau_k) r_{B,B}(\tau_k), \quad (15)$$

where $r_{A,A}(\tau_k)$ is the autocorrelation of the time series of domain A at lag τ_k .

Under the previous null-hypothesis, the expected value of $r_{A,B}(\tau)$ is zero and the following statistic approximately follows the standard normal distribution $N(0, 1)$:

$$z_{A,B}(\tau) = \frac{r_{A,B}(\tau)}{\sqrt{\text{Var}[r_{A,B}(\tau)]}}. \quad (16)$$

The approximation is due to the fact that $r_{A,B}(\tau)$ is bounded between $[-1, 1]$. So, we can now perform hypothesis testing for every pair of domains, computing a corresponding p -value based on z .

Given that there may be several domains in G_δ , we need to control the number of false positive edges that may result from the multiple testing problem. We do so using the False Discovery Rate (FDR) method of Benjamini and Hochberg [19]. Specifically, given N domains, we need to perform $M = \frac{N(N-1)}{2} (2\tau_{max} + 1)$ tests (for each potential edge and for each possible lag value), and compute the p -value for each test, based on (16). Given a False Discovery Rate q (the expected value of the fraction of tests that are false positives), the Benjamini-Hochberg procedure ranks the M p -values (p_i becomes the i 'th lowest p -value) and only keeps the first $m < M$ tests (edges), where p_m is the highest p -value such that $p_m < q m/M$.

Lag inference and edge directionality We infer the domain-level network G_δ as follows. Two domains $A, B \in V_\delta$ are connected if there is at least one lag value at which the cross-correlation $r_{A,B}(\tau)$ has passed the FDR test. The standard approach in *lag inference* is to

consider the lag value τ^* that maximizes the absolute cross-correlation,

$$\tau_{A,B}^* = \arg \max_{\tau=-\tau_{max} \dots \tau_{max}} \{|r_{A,B}(\tau)|\}. \quad (17)$$

The corresponding correlation is denoted as $r_{A,B}^*$. There are two problems with this approach. First, it is harder to examine the statistical significance of $|r_{A,B}^*|$ because it is the maximum of a set of random variables.¹ Second, it is often the case that there is a range of lag values that produce “almost maximum” cross-correlations, say within one standard deviation from each other. Focusing on $\tau_{A,B}^*$ and ignoring the rest of the statistically significant and almost equal cross-correlations is not well justified.

Instead, we follow a more robust approach in which an edge of the domain-level network G_δ may be associated with a range of lag values.² The lag range that we associate with the edge between A and B , denoted as $R_\tau(A, B)$, is defined as *the range of lags that produce significant cross-correlations, within one standard deviation from $|r_{A,B}^*|$* . If $R_\tau(A, B)$ includes $\tau=0$, the edge is represented as *undirected*. If $R_\tau(A, B)$ includes only positive lags, the edge is directed from A to B meaning that A ’s signal precedes B ’s by the given lag range; otherwise, we associate the opposite direction with that edge. We emphasize that the directionality of the edges does *not* imply causality; it only refers to temporal ordering.

Edge weight and domain strength How to assign a weight to each domain-level edge in G_δ ? A common approach is to consider the (signed) magnitude of the cross-correlation $r_{A,B}^*$. This is reasonable if all domain signals have approximately the same signal power. In addition, we propose a new edge weight that is based on the covariance of the two domains:

$$w(A, B) = \text{cov}[X_A(t), X_B(t)] = \tilde{\sigma}_A \tilde{\sigma}_B r_{A,B}^*. \quad (18)$$

¹An analytic approach based on extreme-value statistics was proposed in [100] but it relies on several approximations. Numerical approaches based on frequency-domain bootstrapping, on the other hand, are computationally expensive [100, 107, 127].

²In principle, it may be a set of lag values. In practice though, significant correlations result for a continuous range of lag values.

The cross-correlation is computed at lag $\tau_{A,B}^*$ but we could use the average of all cross-correlations in $R_\tau(A, B)$ instead. The weight of an edge can be positive or negative depending on the sign of the corresponding cross-correlation.

Finally, the strength of a network node (domain) is defined as the sum of the absolute weights of all edges of that node (ignoring edge directionality).

5.4 Illustration - Comparisons

In this section we validate δ -MAPS using the synthetic data set presented in section 2.1. The parameters of δ -MAPS are set as follows: $K=4$ cells (up-down-left-right), and $\delta=0.55$ (corresponds to significance level 10^{-2}). In the edge inference step, the FDR threshold is $q=10\%$ and $\tau_{max} = 20$.

Fig.1-B shows the local homogeneity field $\hat{r}_K(i)$ as well as the identified seeds (blue dots), while Fig.1-C shows the five discovered domains. As expected, we often identify more than one seed in the core of each domain due to noise; those seeds are eventually merged into the same domain. The local homogeneity field is weaker in domains 4 and 5 (due to their lower variance) but a seed is still detected in those domains. Seeds also appear at the two overlapping regions between (1,2) and (2,3) but those seeds gradually merge with one of the domains in which they appear.

Each domain is a subset of the domain's true expanse. The reason is that some cells close to the periphery of each domain have very low signal-to-noise ratio (recall that the signal decays to zero at the periphery and so the average correlation between those cells with the rest of their domain does not exceed the δ threshold). More quantitatively, the inferred domains include about 80%-90% of the ground-truth cells in each domain. In non-overlapping regions this fraction is higher (85%-95% of the cells), while in overlapping regions it drops to 45%-80%. The extent of overlapping regions is harder to correctly identify especially when a domain (e.g., domain 2) overlaps with a stronger domain

(e.g., domains 1 or 3); the stronger domain effectively masks the signal of the weaker domain. The average pairwise cross-correlation of the cells in each domain varies between 55%-70% in the ground-truth data, while the inferred domains have slightly higher average cross-correlation (65%-75%) due to their smaller expanse.

Finally, Fig. 1-C shows the inferred domain-level network. δ -MAPS identifies correctly the three edges and their polarity (positive versus negative correlations). The lag ranges always include the correct value (e.g., the edge between domains 1 and 3 has a lag range [14,15]). Also, the three edges are correctly ordered in terms of absolute cross-correlation magnitude: (1,3) followed by (4,5), followed by (3,5).

5.5 Application in Climate Science

We first apply δ -MAPS in the context of climate science. Climate scientists are interested in *teleconnections* between different regions, and they often rely on EOF analysis to uncover them [167]. Here, we analyze the monthly *Sea-Surface Temperature* (SST) field from the HadISST dataset [121], covering 50 years (1956-2005) at a spatial resolution of $2.0^\circ \times 2.5^\circ$, and we focus on the latitudinal range of $[60^\circ S; 60^\circ N]$ to avoid sea-ice covered regions. Following standard practice, we pre-process the time series to form *anomalies*, i.e., remove the seasonal cycle, remove any long-term trend at each grid-point (using the Theil-Sen estimator), and transform the signal to zero-mean at each grid point.

δ -MAPS is applied as follows. We set the local neighborhood to the $K=4$ nearest cells so that we can identify the smallest possible domains at the given spatial resolution. Second, the homogeneity threshold δ is set to 0.37 (corresponds to a significance level of 10^{-2}). In the edge inference stage, the lag range is $\tau_{max}=12$ months (a reasonable value for large-scale changes in atmospheric wave patterns), and the FDR threshold is set to $q=3\%$ (we identify about 30 edges and so we expect no more than one false positive).

Fig. 57-A shows the identified domains (the color code will be explained shortly). The spatial dimensionality has been reduced from about 6000 grid cells to 18 domains. 65%

of the sea-covered cells belong to at least one domain; the overlapping regions are shown in black and they cover 2% of the grid cells that belong to a domain. The largest domain (domain E) corresponds to the El Niño Southern Oscillation (ENSO), which is also the most important in terms of node strength (see Fig. 57-B). Other strong nodes are domain F (part of the “horseshoe-pattern” surrounding ENSO), domain J (Indian ocean) and domain Q (sub-tropical Atlantic). The strength of the edges associated with ENSO are shown in Fig. 57-C. These observations are consistent with known facts in climate science regarding ENSO and its positive correlation with the Indian ocean and north tropical Atlantic, and negative correlations with the regions that surround it in the Pacific (horseshoe-pattern) [98].

Fig. 57-D shows the inferred domain-level network. The color code represents the (signed) cross-correlation for each edge. The lag range associated with each edge is shown in Fig. 57-E; recall that some edges are not directed because their lag range includes $\tau=0$. The network consists of five weakly-connected components. If we analyze the largest component (which includes ENSO) as a signed network (i.e., some edges are positive and some negative) we see that it is *structurally balanced* [56]. A graph is structurally balanced if it does not contain cycles with an odd number of negative edges.³ A structurally balanced network can be partitioned in a “dipole”, so that positive edges only appear within each pole and negative edges appear only between the two poles. In Fig. 57-A, the nodes of these two poles are colored as blue and green (the smaller disconnected components are shown in other colors).

Focusing on the lag range of each edge, domain Q seems to play a unique role, as it temporally precedes all other domains in the inferred network. Specifically, its activity precedes that of domains D , E and F by about 5-10 months. The lead of south tropical Atlantic SSTs (domain Q) on ENSO has recently received significant attention in climate science [125]. Our results suggest that SST anomalies in domain Q may impact a large

³For instance, if two friends are both enemies with a third person, they form a balanced social triangle.

portion of the climate system.

Switching to lag inference, we say that a triangle is *lag-consistent* if there is at least one value in the lag range associated with each edge that would place the three nodes in a consistent temporal distance with respect to each other. For instance, in the case of the first triangle of Fig. 57-F, the triangle is lag-consistent if the edge from Q to F has a lag of 8 months and the edge between E and F has lag -2 months (meaning that the direction would be from F to E); several other values would make this triangle lag-consistent. We have verified the lag-consistency of every triangle in the climate network. One exception is the triangle between domains (C, D, G) , shown at the bottom of Fig. 57-F. However, the large lag in the edge from C to G can be explained with the triangle between domains (C, E, G) , which is lag-consistent. We emphasize that the temporal ordering that results from these lag relations should not be misinterpreted as causality; we expect that several of the edges we identify are only due to indirect correlations, not associated with a causal interaction between the corresponding two nodes.

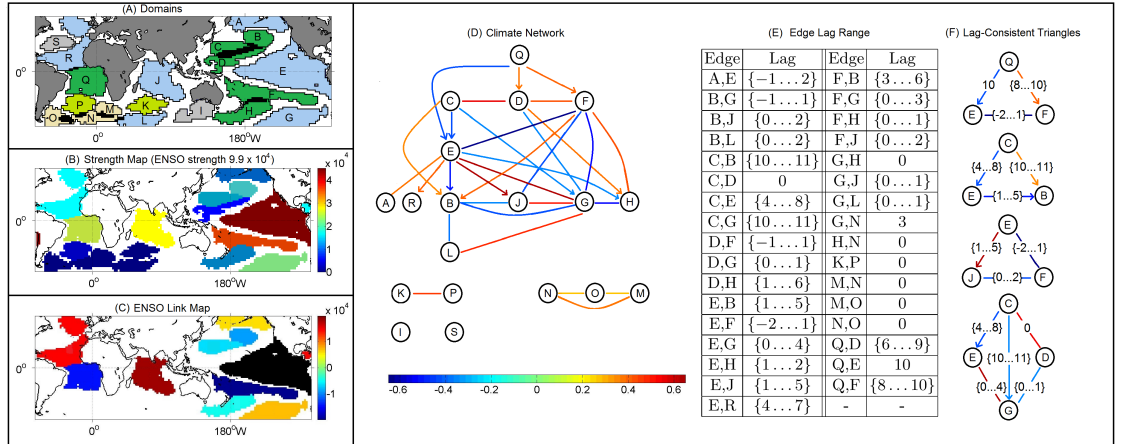


Figure 57: (A) The identified domains. The color of each domain corresponds to the connected component it belongs to (the blue and green nodes belong to two different poles of the same component). (B) Color map for domain strength. The strength of ENSO (domain E) is shown at the top. (C) Edges to and from ENSO (shown in black). (D) The climate network. The color of each edge represents the corresponding cross-correlation. (E) The lag range associated with each edge. (F) Examples of lag-consistent triangles.

For comparison purposes, Fig. 58 shows the results of EOF analysis, community detection, and spatial clustering on the same dataset. The first EOF explains only about 19%

of the variance, implying that the SST field is too complex to be understood with only one spatial component. On the other hand, the joint interpretation of multiple EOF components is problematic due to their orthogonal relation [50]. The anti-correlation between ENSO and the horseshoe-pattern regions is well captured in the first component but several other important connections, such as the negative and lagged relation between the south subtropical Atlantic and ENSO (domains Q and E , respectively), are missed.

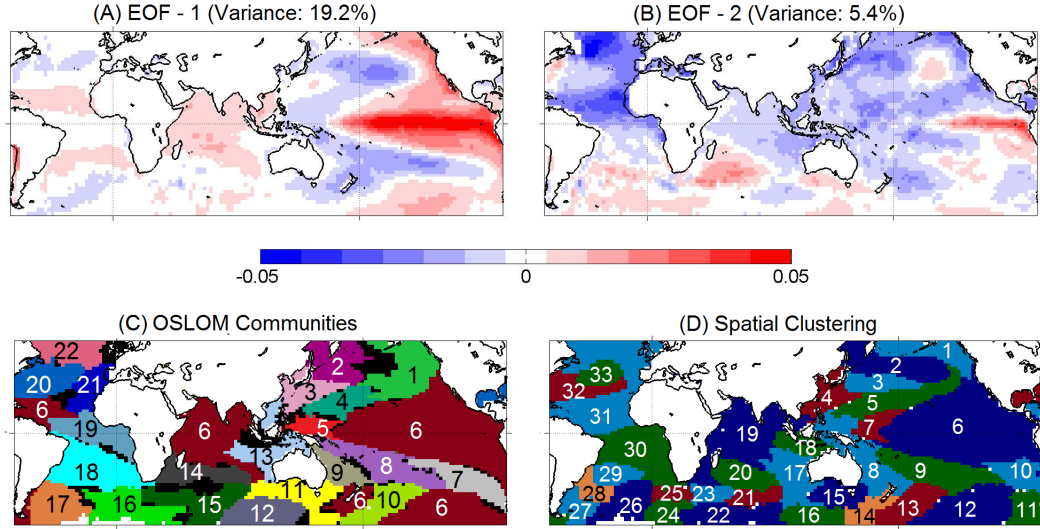


Figure 58: (A),(B) The first two components of EOF analysis. (C) Communities identified by OSLOM. Each community has a unique number and color. (D) Areas identified by spatial clustering.

Fig. 58-C shows the results of the overlapping community detection method OSLOM. Following [142], the input to OSLOM is a correlation-based cell-level network. Correlations less than 30% are ignored. The weight of each edge is set to the maximum absolute correlation between the corresponding two cells, across all considered lags. OSLOM identifies 22 communities. Community 6 is not spatially contiguous; it covers ENSO, the Indian ocean, a region in the north tropical Atlantic, and a region in south Pacific. This is a general problem with community detection methods: they cannot distinguish high correlations due to a remote connection from correlations due to spatial proximity. In the context of climate, the former may be due to atmospheric waves or large-scale currents while the latter may be due to local circulations.

Finally, Fig. 58-D shows the results of a spatial clustering method [66], with the same homogeneity threshold δ we use in δ -MAPS. That method ensures that every cluster (referred to as “area”) is spatially contiguous but it also requires that there is no overlap between areas and it attempts to assign each grid cell to an area. Consequently, it results in more areas (compared to the number of domains), some of which are just artifacts of the spatial parcellation process. Further, the spatial expanse of an area constrains the computation of subsequent areas because no overlaps are allowed.

5.6 Application in fMRI data

Functional magnetic resonance imaging (fMRI) measures fluctuations of the blood oxygenation level dependent (BOLD) signal in the brain. The dynamics of the BOLD signal in gray matter are generally correlated with the level of neural activity. The resulting spatio-temporal field is often analyzed using ICA, clustering or network-based methods to infer *brain functional networks* [136].

Here, we illustrate δ -MAPS on cortical *resting-state* fMRI data from a single subject (healthy young male adult, subject-ID: 122620) from the WU-Minn Human Connectome Project (HCP) [163]. The data acquisition parameters are described in [133]. The spatial resolution is 2mm in each voxel dimension. The pre-processing of fMRI data requires several steps; we use the “fix-extended” HCP minimal processing pipeline that includes head motion correction, registration to a structural image, masking on non-brain voxels, etc; please see [74]. MELODIC ICA and FIX are used to remove non-neuronal artifacts (e.g., physiological noise due to cardiac and respiratory cycles). We also perform bandpass filtering in the range 0.01-0.08Hz, as commonly done in resting-state fMRI.

In this chapter, we analyze two scanning runs of the same subject (“scan-1” and “scan-2”). Each scan lasts about 14 minutes and results in a time series of length $T=1200$ (repetition time $TR=720\text{msec}$). We emphasize that major differences across different scanning sessions of the same subject are common in fMRI; studies of functional brain networks

often only report group-level averages. The entire cortical volume is projected to a surface mesh (Conte69 32K) resulting in about 65K *gray-ordinate* points (as opposed to volumetric voxels) [162]. Each point of this mesh is adjacent to six other points; for this reason we set $K=6$. The homogeneity threshold is set to $\delta=0.37$ (corresponds to significance level 10^{-2}). The maximum lag range τ_{max} is set to ± 3 , i.e., 2.2 seconds, and the FDR threshold is set to $q=10^{-4}$ (i.e., we expect one out of 10K edges to be a false positive). The signal of a domain is defined as the average across all voxels in that domain.

The application of δ -MAPS results in a network with about 850 domains in scan-1 (1120 domains in scan-2). 80% of the domains are smaller than 30-40 voxels (depending on the scan) and 5% of the domains are larger than 250 voxels. The number of edges is 4285 in scan-1 (4200 in scan-2). The absolute value of the cross-correlation associated with each edge is typically larger than 0.5. The fraction of negative edge correlations is about 5% in scan-1 and 20% in scan-2 suggesting that the polarity of some network edges may be time-varying. The lag τ^* that corresponds to the maximum cross-correlation is 0 in 70% of the edges and ± 1 in almost all other cases. 13% of the edges are directed, meaning that lag-0 does not produce a significant correlation for that pair of domains. There is a positive correlation between the degree of a domain and its physical size (the correlation coefficient between degree and $\log_{10}(\text{size})$ is 0.70 for scan-1 and 0.66 for scan-2). Further, the network is assortative meaning that domains tend to connect to other domains of similar degree (assortativity coefficient about 0.7 in both scans).

An important question is whether the δ -MAPS networks are consistent with what neuroscientists currently know about resting-state activity in the brain. During rest, certain cortical regions that are collectively referred to as the *Default-Mode Network (or DMN)* are persistently active across age and gender [176]. Other known resting-state networks are the occipital (part of the visual system) and the motor/somatosensory (associated with planning and execution of voluntary body motion). With the terminology of network theory, the previous “networks” would be referred to as *communities* within the larger functional brain

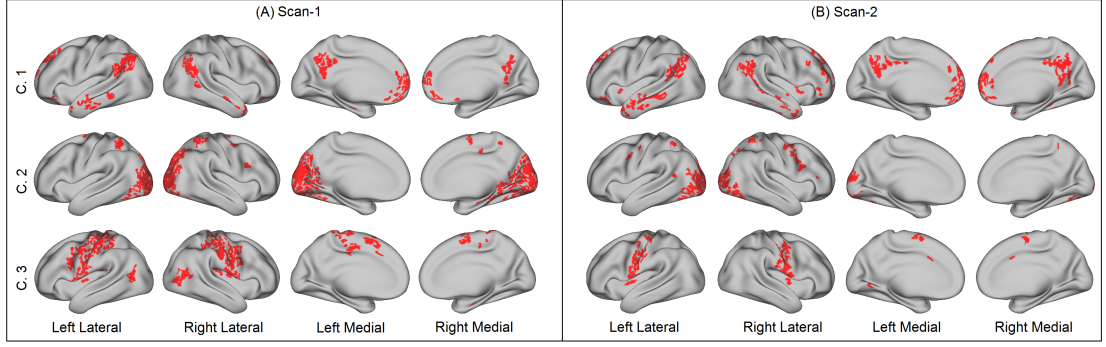


Figure 59: Three domain-level network communities for each scan. The first corresponds to the default-mode network, the second to the occipital network, and the third to the motor/somatosensory network.

network. To identify communities in the δ -MAPS network, we applied OSLOM [103]. OSLOM identifies two hierarchical levels in both scans. The first level consists of highly overlapping communities that cover almost the entire cortex. The second hierarchical level is more interesting, resulting in eight communities for scan-1 (nine for scan-2). Fig. 59 shows the three communities (C.1, C.2, C.3) for each scan that have the highest resemblance to the three previously mentioned resting-state networks: C.1 corresponds to the DMN, C.2 corresponds to the occipital resting-state network, and C.3 corresponds to the motor/somatosensory network. C.1 is quite similar across the two scanning sessions and it clearly captures the DMN. In C.2, the extent of the network is smaller in scan-2, which is not too surprising giving the known inter-scan variability of resting-state fMRI. C.3 is also quite similar across the two scans and consistent with the motor/somatosensory network.

To further investigate the structure of those higher degree (and typically larger) domains, we perform *k-core decomposition*.⁴ The density of the remaining network, after the extraction of $k=14$ cores from the scan-1 network ($k=16$ cores in scan-2) shows a sudden increase by a factor of two. This suggests that the network includes a *densely inter-connected backbone*, also known as “rich-club”. The size of this backbone is small relative to the entire network: 130 domains in scan-1 (90 in scan-2). Similar observations about the resting-state brain, but using voxel-level network analysis methods, have been

⁴A process that starts with the original network ($k=0$), and it removes iteratively all nodes of degree k or less in each round so that after the extraction of the k 'th core all remaining nodes have degree larger than k .

previously reported [161]. Fig.60 shows the location of the backbone domains for each hemisphere and for each scan. The regions that are usually associated with the DMN dominate the backbone of both sessions. Interestingly though, scan-1 includes the regions of the motor/somatosensory network, while the backbone of scan-2 is missing those regions. One possible explanation for this discrepancy is that the subject was more relaxed during scan-2, not exerting the mental effort to stay still.

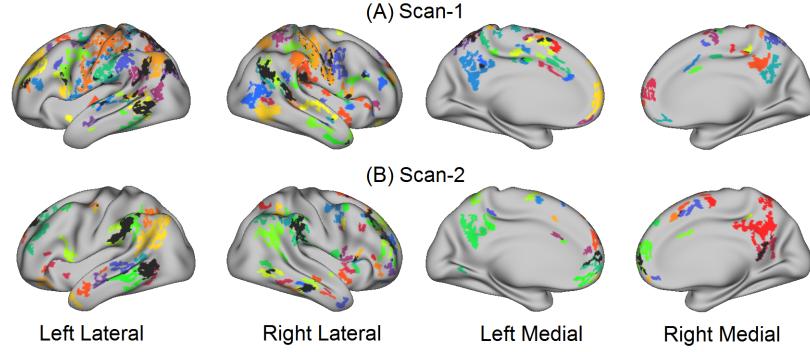


Figure 60: The domains of the backbone network for each hemisphere and scan. The color of each domain is randomly assigned (overlaps are shown in black).

5.7 Discussion

δ -MAPS results in a correlation-based functional network. A next step could be to infer a causal, or *effective* network, leveraging the framework of probabilistic graphical models. Instead of attempting to learn the graph structure from raw data, one could use the δ -MAPS network as the underlying structure and then apply conditional independence tests to remove non-causal edges (e.g., [58]). Another direction could be to combine the inferred functional network with a structural network that shows the physical connectivity between the identified domains. This is not hard in the case of communication networks but it also becomes feasible for brain networks using diffusion-weighted MRI. The projection of the observed dynamics on the underlying structure can help to characterize the actual function and delay of each system component.

5.8 Identifying the largest domain is NP-complete

We are given a spatio-temporal field $\mathbf{X}(t)$ on a grid G , a pairwise similarity metric between pairs of grid cells and a threshold δ . Starting from a grid cell c , the goal is to find the largest subset of grid cells that form a single spatially connected component, and whose average similarity exceeds the threshold δ . The spatial grid can be represented as a planar graph $G(V, E)$ where each grid cell is a node and edges connect adjacent grid cells. Formally we have the following graph optimization problem:

Definition 1. Rooted Largest Connected δ -Dense Subgraph Problem (rooted LC δ DS).

Given a regular (grid) graph $G(V, E)$, a weight function $w : V \times V \rightarrow \mathbb{R}$ (where $w(v, v) = 0$ and symmetric), a threshold δ , and a node $c \in V$, find a maximum cardinality set of nodes $A \subseteq V$ such that $c \in A$, the induced subgraph is connected ($I_G(A) = 1$) and $\frac{\sum_{v,u \in A} w(v,u)}{|A|(|A|-1)} > \delta$ (i.e., $\hat{r}(A) > \delta$).

To show that rooted LC δ DS is NP-hard we first consider a variant of the problem in which the induced subgraph A has to satisfy two conditions; it has to be a connected subgraph of G , and the average weight of the edges in A has to exceed δ . More formally:

Definition 2. Largest Connected δ -Dense Subgraph Problem (LC δ DS). Given a regular (grid) graph $G(V, E)$, a weight function $w : V \times V \rightarrow \mathbb{R}$ (where $w(v, v) = 0$ and symmetric), and a threshold δ , find a maximum cardinality set of nodes $A \subseteq V$ such that $I_G(A) = 1$ and $\hat{r}(A) > \delta$.

To show that LC δ DS is NP-hard we use a reduction of the densest connected k subgraph problem.

Definition 3. Densest Connected k -Subgraph Problem (DC k S). Decision version: Given a graph $G(V, E)$, and positive integers k and j , does there exist an induced subgraph on k vertices such that this subgraph has at least j edges and is connected?

DC k S (also referred to as the connected h-clustering problem) has been shown to be NP-complete on general graphs [42], as well as on planar graphs [96]. DC k S is polynomially time solvable for subclasses of planar graphs of bounded tree width [12]. Grid

graphs, which are the type of graphs that arise in our application domains, are planar bipartite graphs, with non-fixed tree width, and no positive results are known for this subclass of planar graphs. The work on approximating densest/heaviest connected k -subgraphs is relatively very limited (see recent theoretical result [36]). It is easy to show that the $DCkS$ problem can be easily reduced to an instance of the decision version of the $LC\delta DS$ problem, and hence it is also NP-complete even on planar graphs.

LEMMA 1. The decision version of the $LC\delta DS$ problem is NP-complete on planar graphs.

PROOF. This can be shown via a reduction from the $DCkS$. We reduce an instance $\langle G, k, j \rangle$ of the $DCkS$ to an $LC\delta DS$ instance by using the same graph G , setting $w(u, v) = I(u, v) \in E$ ($w(u, v)$ is 1 if and only if the pair of nodes is connected by an edge), and $\delta = j/k(k-1)$.

Now it is easy to show that rooted $LC\delta DS$ is also NP-hard. If a poly-time algorithm existed for the rooted $LC\delta DS$, then by calling it $|V|$ times with each of the nodes of the graph, we would obtain in poly-time a solution to the NP-hard $LC\delta DS$.

5.9 Heuristic for the selection of δ

The threshold δ intuitively determines the minimum degree of homogeneity that the underlying field must have within each domain. The higher the threshold, the higher the required homogeneity and therefore, the smaller the size of the identified domains.

To select δ we propose the following heuristic. We start with a random sample of pairs of grid cells and for each pair i, j we compute the Pearson correlation $r_{i,j}$ at zero lag. To assess the significance of each correlation we use Bartlett's formula [26]. Under the null hypothesis of no coupling $r_{i,j}$ should have zero mean, and a reasonable estimate of its variance is given by

$$Var[r_{i,j}] = \frac{1}{T} \sum_{\tau_k=-T}^T r_{i,i}(\tau_k) r_{j,j}(\tau_k), \quad (19)$$

here $r_{i,i}(\tau_k)$ is the autocorrelation of the time series of grid cell i at lag τ_k . The scaled values $z_{i,j} = \frac{r_{i,j}}{\sqrt{\text{Var}[r_{i,j}]}}$ should approximately follow a standard normal distribution. To assess the significance of each correlation we perform a one sided z-test for a given level of significance α .

The threshold δ is set as the average of all significant correlations. A domain is a set of spatially contiguous grid cells, thus we require that the mean pairwise correlation for the cells belonging to the same domain to be higher than the mean pair-wise correlation of randomly picked pairs of grid cells. δ depends on the choice of the significance level α , on the autocorrelation structure of the underlying time series and on the correlation distribution of the field.

5.10 δ -MAPS pseudocode

```

1: Domains  $S = \{A_1, \dots, A_{|S|}\}$  ▷ The initial set of domains
2: function DOMAINIDENTIFICATION()
3:   while True do
4:     boolean  $merged \leftarrow$  DOMAINMERGING( $S$ )
5:     boolean  $expanded \leftarrow$  DOMAINEXPANSION( $S$ )
6:     if  $!merged \&\& !expanded$  then
7:       break ▷ Terminate when no further expansion or merging is
           possible
8:     end if
9:   end while
10: end function

```

```

1: function DOMAINEXPANSION(Domains  $S = \{A_1, \dots, A_{|S|}\}$ )
2:   boolean startMerging  $\leftarrow false$ 
3:   boolean expanded  $\leftarrow false$ 
4:   while !startMerging do  $\triangleright$  Domain expansion is repeated in rounds
5:     expanded  $\leftarrow false$ 
6:     sort( $S$ )  $\triangleright$  Sort domains in decreasing order of homogeneity such
       that  $\hat{r}(A_{i-1}) > \hat{r}(A_i) > \hat{r}(A_{i+1})$ 
7:     for  $i = 1 : |S|$  do
8:       Domain  $A_i \leftarrow S[i]$ 
9:       Domain  $eA_i \leftarrow \text{EXPANDDOMAIN}(A_i)$ 
10:      if  $|A_i| \neq |eA_i|$  then  $\triangleright$  Domain expanded
11:         $S[i] \leftarrow eA_i$ 
12:        expanded  $\leftarrow true$ 
13:        startMerging  $\leftarrow \text{CANMERGE}(eA_i)$ 
14:        if startMerging then
15:          break  $\triangleright$  Exit the for loop
16:        end if
17:      end if
18:    end for  $\triangleright$  A round of domain expansion is complete
19:    if !expanded then
20:      break  $\triangleright$  Domains cannot be expanded
21:    end if
22:  end while
23:  return expanded
24: end function
25:
26: function EXPANDDOMAIN(Domain  $A_i$ )  $\triangleright$  Try to expand domain  $A_i$  by
    one cell
27:   Construct set  $\Gamma(A_i)$ : all cells adjacent to  $A_i$ 
28:   if  $\Gamma(A_i) = \emptyset$  then
29:     return  $A_i$ 
30:   else
31:      $m \leftarrow \arg \max_{m \in \Gamma(A_i)} \hat{r}(A_i \cup \{m\})$   $\triangleright$  Select the cell that maximizes
        $\hat{r}(A_i \cup \{m\})$ .
32:     if  $\hat{r}(A_i \cup \{m\}) > \delta$  then
33:        $A_i \leftarrow A_i \cup m$ 
34:     end if
35:     return  $A_i$ 
36:   end if
37: end function
38:
39: function CANMERGE(Domain  $A_i$ )  $\triangleright$  Check whether one or more merging
    operations are possible
40:   boolean merge  $\leftarrow false$ 
41:   Construct set  $\Gamma(A_i)$ : all domains adjacent to  $A_i$ 
42:   for  $j = 1 : |\Gamma(A_i)|$  do
43:      $A_j \leftarrow \Gamma(A_i)[j]$ 
44:     if  $\hat{r}(A_i \cup A_j) > \delta$  then 2
45:       merge  $\leftarrow true$ 
46:       break
47:     end if
48:   end for
49:   return merge
50: end function

```

```

1: function DOMAINMERGING(Domains  $S = \{A_1, \dots, A_{|S|}\}$ )
2:   boolean  $merged \leftarrow false$ 
3:   while True do  $\triangleright$  Repeat until no pair of domains can be merged
4:     Domain  $DomainToMerge1 \leftarrow \emptyset$ 
5:     Domain  $DomainToMerge2 \leftarrow \emptyset$   $\triangleright$  Domains with the maximum
      union homogeneity
6:      $maxHomogeneity \leftarrow -1$ 
7:     for  $i = 1 : |S|$  do
8:       Domain  $A_i \leftarrow S[i]$   $\triangleright$  Get the  $i^{th}$  domain
9:       Construct set  $\Gamma(A_i)$ 
10:       $A_j \leftarrow \arg \max_{A_j \in \Gamma(A_i)} \hat{r}(A_i \cup A_j)$ 
11:      if  $\hat{r}(A_i \cup A_j) > maxHomogeneity$  then  $\triangleright$  Update the best
        candidates to merge
12:         $DomainToMerge1 \leftarrow A_i$ 
13:         $DomainToMerge2 \leftarrow A_j$ 
14:         $maxHomogeneity \leftarrow \hat{r}(A_i \cup A_j)$ 
15:      end if
16:    end for
17:    if  $maxHomogeneity > \delta$  then
18:      S.remove( $DomainToMerge1$ )
19:      S.remove( $DomainToMerge2$ )  $\triangleright$  Remove the domains that will
        be merged
20:       $S \leftarrow DomainToMerge1 \cup DomainToMerge2$ 
21:       $merged \leftarrow true$ 
22:    else
23:      break  $\triangleright$  We can not merge any domains
24:    end if
25:  end while
26:  return  $merged$   $\triangleright$  Return true if at least one pair of domains is merged
27: end function

```

Chapter VI

CONCLUSIONS & FUTURE WORK

6.1 Conclusions

In this thesis we propose a framework for the analysis of spatio-temporal systems based on complex network analysis. The proposed framework consists of two methods *geo-Cluster* and δ -MAPS, whose scope is to uncover the semi-autonomous functional components of a spatio-temporal system and infer their interactions.

The first method, *geo-Cluster*, identifies the functional components of the system, referred to as “areas”, and models their interconnections as a complete and weighted network. An area is a spatially contiguous, non-overlapping, set of grid cells that conform to a homogeneity constraint. This homogeneity constraint requires that the average pairwise correlation between the grid cells in an area’s scope to be larger than a pre-defined threshold - the only parameter of the proposed algorithm. The requirement of only one parameter, combined with the fact that no link pruning in the underlying cell-level network is imposed, adds robustness to a network’s structure and makes the comparison of different networks more reliable. At a second step, we infer a network between these areas. The network is modeled as a complete and weighted graph. The weight of an edge, measured as the covariance between the time series of the two corresponding areas, captures the magnitude of the interaction between the functional components of the system.

The proposed method is robust to noise, the resolution of the spatio-temporal data set, the measure that quantifies similarities between the grid cell time series, and to perturbations of the homogeneity parameter.

The second method, δ -MAPS, allows for the functional components of the system (referred to as “domains”) to overlap and accounts for non-instantaneous interactions between

them. δ -MAPS is based on the premise that the functional relation between the grid cells of a domain results in highly correlated temporal activity. To this end it first identifies the “epicenter” or “core” of a domain as a point (or set of points) where the local homogeneity is maximum across the entire domain. Instead of searching for the discrete boundary of a domain, which may not exist in reality, we compute a domain as the *maximum possible set* of spatially contiguous cells that include the detected core, and that satisfy a homogeneity constraint, expressed in terms of the average pairwise cross-correlation across the domain’s scope. At a second step, δ -MAPS infers a functional network. Different domains may have correlated activity, potentially at a lag, because of direct or indirect interactions. The proposed edge inference method examines the statistical significance of each lagged cross-correlation between two domains, applies a multiple-testing process to control the rate of false positives, infers a range of potential lag values for each edge, and assigns a weight to each edge based on the covariance of the corresponding two domains.

δ -MAPS does not require the number of domains as an input parameter, the resulting domains are spatially contiguous and potentially overlapping, and the inferred connections between domains can be lagged and positively or negatively weighted. Further, the distinction between grid cells that are correlated within the same domain and grid cells that are correlated across two distinct domains allows δ -MAPS to separate between local diffusion (or dispersion) phenomena and remote interactions that may be due to underlying structural connections (e.g., a white-matter fiber between two brain regions).

δ -MAPS is not just a generalization of *geo-Cluster*, allowing for overlapping functional domains and accounting for lagged interactions between them. The greedy heuristics of *geo-Cluster* force each grid cell to belong to an area. Further, after a grid cell is assigned to an area it cannot belong to any other area, potentially limiting the scope of subsequent areas. This leads to a stronger path dependency (thus less robustness) compared to the approach taken by the δ -MAPS algorithm.

The proposed framework has been applied in the fields of climate science and neuroscience. In the context of climate we present applications of *geo-Cluster* to identify well known climate shifts and construct networks between different climate fields. Using *geo-Cluster* we performed an extensive study analyzing twelve cutting edge climate model ensembles from the CMIP5 output. Using two distance metrics, a network distance D and the adjusted Rand index (ARI) we are able to rank the models in terms of their ability to reproduce the climate of the past as well as quantify the variability between different members of the same model ensemble. When investigating the model trajectories in the future, under a global warming scenario, we found that the uncertainty in the model trajectories is larger than the uncertainty in the superimposed trends.

Using δ -MAPS we analyzed the temporal relationships between different functional components of the climate system in the sea surface temperature field. We found that the proposed method successfully uncovered many well-known climate teleconnections and the lag associated with them. In the context of neuroscience we performed a single subject analysis focusing on resting state fMRI data. We found that the proposed method was able to uncover many of the well-known resting state networks. We also show how the method identifies a small number of strongly interconnected areas forming the backbone of the resting state network. Using synthetic data we also show how δ -MAPS overcomes limitations of traditional dimensionality reduction techniques such as PCA/ICA, clustering and community detection.

6.2 Future Work

Climate Networks Over Time. The proposed method can be naturally extended to construct networks over time (e.g., using a sliding window approach). To this end, we are interested to observe the trajectories of the functional components of the system as expressed by their strength and size. Specifically, in the context of climate we are interested to focus on the dynamics of ENSO and identify “tipping” points at which its dynamics

change. Such tipping points can have a global impact on the climate system.

Climate Models and Controlled Perturbation Experiments. In the context of climate we are interested to use the network framework to evaluate how the perturbations imposed on a model’s parameters propagate to the climate scale. We are interested to first identify the regions that are the most (or the least) affected by the perturbations, the time scale of the propagation, and the implications for teleconnections.

Effective Connectivity. In this thesis we have limited our analysis to functional networks. A next step could be to infer a causal, or effective [130] network, leveraging the framework of probabilistic graphical models [57, 101, 128]. Instead of attempting to learn the graph structure from raw data, one could use the identified spatial components as the underlying structure and then apply conditional independence tests to remove non-causal edges.

Dynamic Networks Using Contextual Time Series Detection. A problem that arises with fMRI measurements is that they are sampled over an extended recording period. Most fMRI studies require that the subject will remain at rest through this period (which cannot be guaranteed in practice). When we measure cross-correlations throughout that extended measurement period, abrupt changes in the signal might be averaged out [151]. A proposed solution to this problem, to be able to identify these dynamic changes, is to construct temporal networks using a sliding window approach [33]. However, the results are highly dependent on the length of the window chosen. An alternative direction would be to automatically detect changes between two time series [34, 35]. Such changes can be tracked at the voxel level. However, due to the high amount of noise in fMRI data, a better approach would be to track such changes at the functional domain level.

Structural-Functional Networks. Another direction could be to combine the inferred functional network with a structural network that shows the physical connectivity between the identified domains. This is not hard in the case of communication networks but it also becomes feasible for brain networks using diffusion-weighted MRI. The projection of the

observed dynamics on the underlying structure can help to characterize the actual function and delay of each system component.

Extensions to Other Spatio-temporal Data. The applications of the proposed framework are not only limited in the fields of climate science and neuroscience. To this end, we propose to apply the proposed framework to data describing species migration patterns (see e.g., [91]). By understanding the processes that drive such patterns we can mitigate risks to populations due to climate change, urban expansion and many more factors. In such a context, the functional components that we identify will correspond to migratory regions. The edges between the identified regions can uncover pathways of population movement.

REFERENCES

- [1] ABRAMOV, R. V. and MAJDA, A. J., “A new algorithm for low-frequency climate response,” *Journal of the Atmospheric Sciences*, vol. 66, no. 2, pp. 286–309, 2009.
- [2] ADAMS, R. and BISCHOF, L., “Seeded region growing,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 16, no. 6, pp. 641–647, 1994.
- [3] AHLGRIMM, M. and FORBES, R., “The impact of low clouds on surface shortwave radiation in the ecmwf model,” *Monthly Weather Review*, vol. 140, no. 11, pp. 3783–3794, 2012.
- [4] AHN, Y.-Y., BAGROW, J. P., and LEHMANN, S., “Link communities reveal multi-scale complexity in networks,” *Nature*, vol. 466, no. 7307, pp. 761–764, 2010.
- [5] AKHSHABI, S. and DOVROLIS, C., “The evolution of layered protocol stacks leads to an hourglass-shaped architecture,” in *Dynamics On and Of Complex Networks, Volume 2*, pp. 55–88, Springer, 2013.
- [6] AKRITAS, M. G., MURPHY, S. A., and LAVALLEY, M. P., “The Theil-Sen estimator with doubly censored data and applications to astronomy,” *Journal of the American Statistical Association*, vol. 90, no. 429, pp. 170–177, 1995.
- [7] ALBERT, R. and BARABÁSI, A.-L., “Statistical mechanics of complex networks,” *Reviews of modern physics*, vol. 74, no. 1, p. 47, 2002.
- [8] ALEXANDER-BLOCH, A., LAMBIOTTE, R., ROBERTS, B., GIEDD, J., GOGTAY, N., and BULLMORE, E., “The discovery of population differences in network community structure: new methods and applications to brain functional networks in schizophrenia,” *Neuroimage*, vol. 59, no. 4, pp. 3889–3900, 2012.
- [9] ALLAN, R., LINDESAY, J., PARKER, D., and OTHERS, *El Niño southern oscillation & climatic variability*. CSIRO publishing, 1996.
- [10] ALLEN, M. R. and SMITH, L. A., “Investigating the origins and significance of low-frequency modes of climate variability,” *Geophysical Research Letters*, vol. 21, no. 10, pp. 883–886, 1994.
- [11] ANDRONOVA, N. G. and SCHLESINGER, M. E., “Objective estimation of the probability density function for climate sensitivity,” *Journal of Geophysical Research: Atmospheres*, vol. 106, no. D19, pp. 22605–22611, 2001.
- [12] ARNBORG, S., LAGERGREN, J., and SEESE, D., “Easy problems for tree-decomposable graphs,” *Journal of Algorithms*, vol. 12, no. 2, pp. 308–340, 1991.

- [13] ARTHUR, D. and VASSILVITSKII, S., “k-means++: The advantages of careful seeding,” in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1027–1035, Society for Industrial and Applied Mathematics, 2007.
- [14] BALDASSANO, C., BECK, D. M., and FEI-FEI, L., “Parcellating connectivity in spatial maps,” *PeerJ*, vol. 3, p. e784, 2015.
- [15] BANDUKWALA, F., “Extracting spatially and spectrally coherent regions from multispectral images,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on*, pp. 82–87, IEEE, 2011.
- [16] BARTHÉLEMY, M., “Spatial networks,” *Physics Reports*, vol. 499, no. 1, pp. 1–101, 2011.
- [17] BELLEC, P., PERLBARG, V., JBABDI, S., PÉLÉGRINI-ISSAC, M., ANTON, J.-L., DOYON, J., and BENALI, H., “Identification of large-scale networks in the brain using fmri,” *Neuroimage*, vol. 29, no. 4, pp. 1231–1243, 2006.
- [18] BELLENGER, H., GUILYARDI, É., LELOUP, J., LENGAGNE, M., and VIALARD, J., “Enso representation in climate models: from cmip3 to cmip5,” *Climate Dynamics*, vol. 42, no. 7-8, pp. 1999–2018, 2014.
- [19] BENJAMINI, Y. and HOCHBERG, Y., “Controlling the false discovery rate: A practical and powerful approach to multiple testing,” *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 289–300, 1995.
- [20] BEREZIN, Y., GOZOLCHIANI, A., GUEZ, O., and HAVLIN, S., “Stability of climate networks with time,” *Scientific reports*, vol. 2, 2012.
- [21] BETZEL, R. F., BYRGE, L., HE, Y., GOÑI, J., ZUO, X.-N., and SPORNS, O., “Changes in structural and functional connectivity among resting-state networks across the human lifespan,” *NeuroImage*, vol. 102, pp. 345–357, 2014.
- [22] BIRANT, D. and KUT, A., “ST-DBSCAN: An algorithm for clustering spatial-temporal data,” *Data & Knowledge Engineering*, vol. 60, no. 1, pp. 208–221, 2007.
- [23] BLUMENSATH, T., BEHRENS, T. E., and SMITH, S. M., “Resting-state fmri single subject cortical parcellation based on region growing,” in *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2012*, pp. 188–195, Springer, 2012.
- [24] BLUMENSATH, T., JBABDI, S., GLASSER, M. F., VAN ESSEN, D. C., UGURBIL, K., BEHRENS, T. E., and SMITH, S. M., “Spatially constrained hierarchical parcellation of the brain with resting-state fmri,” *Neuroimage*, vol. 76, pp. 313–324, 2013.
- [25] BOERS, N., BOOKHAGEN, B., MARWAN, N., KURTHS, J., and MARENGO, J., “Complex networks identify spatial patterns of extreme rainfall events of the south american monsoon system,” *Geophysical Research Letters*, vol. 40, no. 16, pp. 4386–4392, 2013.

- [26] BOX, G. E., JENKINS, G. M., and REINSEL, G. C., *Time series analysis: forecasting and control*, vol. 734. John Wiley & Sons, 2011.
- [27] BRACCO, A., KUCHARSKI, F., MOLTENI, F., HAZELEGER, W., and SEVERIJNS, C., “Internal and forced modes of variability in the indian ocean,” *Geophysical research letters*, vol. 32, no. 12, 2005.
- [28] BULLMORE, E. and SPORNS, O., “Complex brain networks: graph theoretical analysis of structural and functional systems,” *Nature Reviews Neuroscience*, vol. 10, no. 3, pp. 186–198, 2009.
- [29] CAI, W., BORLACE, S., LENGAGNE, M., VAN RENSCH, P., COLLINS, M., VECCHI, G., TIMMERMANN, A., SANTOSO, A., MCPHADEN, M. J., WU, L., and OTHERS, “Increasing frequency of extreme el niño events due to greenhouse warming,” *Nature Climate Change*, vol. 4, no. 2, pp. 111–116, 2014.
- [30] CARRÉ, M., SACHS, J. P., PURCA, S., SCHAUER, A. J., BRACONNOT, P., FALCÓN, R. A., JULIEN, M., and LAVALLÉE, D., “Holocene history of enso variance and asymmetry in the eastern tropical pacific,” *Science*, vol. 345, no. 6200, pp. 1045–1048, 2014.
- [31] CARTON, J. A. and GIESE, B. S., “A reanalysis of ocean climate using simple ocean data assimilation (soda),” *Monthly Weather Review*, vol. 136, no. 8, pp. 2999–3017, 2008.
- [32] CHAMBERS, D., TAPLEY, B., and STEWART, R., “Anomalous warming in the indian ocean coincident with el nino,” *Journal of Geophysical Research: Oceans*, vol. 104, no. C2, pp. 3035–3047, 1999.
- [33] CHANG, C. and GLOVER, G. H., “Time–frequency dynamics of resting-state brain connectivity measured with fmri,” *Neuroimage*, vol. 50, no. 1, pp. 81–98, 2010.
- [34] CHEN, X. C., MUEEN, A., NARAYANAN, V. K., KARAMPATZIAKIS, N., BANSAL, G., and KUMAR, V., “Online discovery of group level events in time series,” in *SDM*, pp. 632–640, SIAM, 2014.
- [35] CHEN, X. C., STEINHAEUSER, K., BORIAH, S., CHATTERJEE, S., and KUMAR, V., “Contextual time series change detection,” in *SDM*, pp. 503–511, SIAM, 2013.
- [36] CHEN, X., HU, X., and WANG, C., “Finding connected dense k-subgraphs,” in *Theory and Applications of Models of Computation*, pp. 248–259, Springer, 2015.
- [37] COBB, K. M., WESTPHAL, N., SAYANI, H. R., WATSON, J. T., DI LORENZO, E., CHENG, H., EDWARDS, R., and CHARLES, C. D., “Highly variable el niño–southern oscillation throughout the holocene,” *Science*, vol. 339, no. 6115, pp. 67–70, 2013.

- [38] COHEN, A. L., FAIR, D. A., DOSENBACH, N. U., MIEZIN, F. M., DIERKER, D., VAN ESSEN, D. C., SCHLAGGAR, B. L., and PETERSEN, S. E., “Defining functional areas in individual human brains using resting functional connectivity mri,” *Neuroimage*, vol. 41, no. 1, pp. 45–57, 2008.
- [39] COLLINS, M., AN, S.-I., CAI, W., GANACHAUD, A., GUILYARDI, E., JIN, F.-F., JOCHUM, M., LENGAGNE, M., POWER, S., TIMMERMANN, A., and OTHERS, “The impact of global warming on the tropical pacific ocean and el niño,” *Nature Geoscience*, vol. 3, no. 6, pp. 391–397, 2010.
- [40] COLLINS, M. and OTHERS, “El niño-or la niña-like climate change?,” *Climate Dynamics*, vol. 24, no. 1, pp. 89–104, 2005.
- [41] CORMEN, T. H., LEISERSON, C. E., RIVEST, R. L., and STEIN, C., *Introduction to algorithms*, vol. 6. MIT press Cambridge, 2001.
- [42] CORNEIL, D. G. and PERL, Y., “Clustering and domination in perfect graphs,” *Discrete Applied Mathematics*, vol. 9, no. 1, pp. 27–39, 1984.
- [43] CORTI, S., GIANNINI, A., TIBALDI, S., and MOLTENI, F., “Patterns of low-frequency variability in a three-level quasi-geostrophic model,” *Climate Dynamics*, vol. 13, no. 12, pp. 883–904, 1997.
- [44] CRADDOCK, R. C., JAMES, G. A., HOLTZHEIMER, P. E., HU, X. P., and MAYBERG, H. S., “A whole brain fmri atlas generated via spatially constrained spectral clustering,” *Human brain mapping*, vol. 33, no. 8, pp. 1914–1928, 2012.
- [45] CROSSLEY, N. A., MECHELLI, A., SCOTT, J., CARLETTI, F., FOX, P. T., MCGUIRE, P., and BULLMORE, E. T., “The hubs of the human connectome are generally implicated in the anatomy of brain disorders,” *Brain*, vol. 137, no. 8, pp. 2382–2395, 2014.
- [46] DEE, D., UPPALA, S., SIMMONS, A., BERRISFORD, P., POLI, P., KOBAYASHI, S., ANDRAE, U., BALMASEDA, M., BALSAMO, G., BAUER, P., and OTHERS, “The era-interim reanalysis: Configuration and performance of the data assimilation system,” *Quarterly Journal of the Royal Meteorological Society*, vol. 137, no. 656, pp. 553–597, 2011.
- [47] DENG, Y. and EBERT-UPHOFF, I., “Weakening of atmospheric information flow in a warming climate in the community climate system model,” *Geophysical Research Letters*, vol. 41, no. 1, pp. 193–200, 2014.
- [48] DESER, C., PHILLIPS, A. S., TOMAS, R. A., OKUMURA, Y. M., ALEXANDER, M. A., CAPOTONDI, A., SCOTT, J. D., KWON, Y.-O., and OHBA, M., “Enso and pacific decadal variability in the community climate system model version 4,” *Journal of Climate*, vol. 25, no. 8, pp. 2622–2651, 2012.

- [49] DIJKSTRA, H. A., *Nonlinear physical oceanography: a dynamical systems approach to the large scale ocean circulation and El Nino*, vol. 28. Springer Science & Business Media, 2005.
- [50] DOMMENGET, D. and LATIF, M., “A cautionary note on the interpretation of EOFs,” *Journal of Climate*, vol. 15, no. 2, pp. 216–225, 2002.
- [51] DONGES, J. F., SCHULTZ, H. C., MARWAN, N., ZOU, Y., and KURTHS, J., “Investigating the topology of interacting networks,” *The European Physical Journal B*, vol. 84, no. 4, pp. 635–651, 2011.
- [52] DONGES, J. F., ZOU, Y., MARWAN, N., and KURTHS, J., “The backbone of the climate network,” *EPL (Europhysics Letters)*, vol. 87, no. 4, p. 48007, 2009.
- [53] DONGES, J. F., ZOU, Y., MARWAN, N., and KURTHS, J., “Complex networks in climate dynamics,” *The European Physical Journal Special Topics*, vol. 174, no. 1, pp. 157–179, 2009.
- [54] DOWNAR, J., CRAWLEY, A. P., MIKULIS, D. J., and DAVIS, K. D., “A multimodal cortical network for the detection of changes in the sensory environment,” *Nature neuroscience*, vol. 3, no. 3, pp. 277–283, 2000.
- [55] DUQUE, J. C., RAMOS, R., and SURINACH, J., “Supervised regionalization methods: A survey,” *International Regional Science Review*, vol. 30, no. 3, pp. 195–220, 2007.
- [56] EASLEY, D. and KLEINBERG, J., *Networks, crowds, and markets: Reasoning about a highly connected world*. Cambridge University Press, 2010.
- [57] EBERT-UPHOFF, I. and DENG, Y., “Causal discovery for climate research using graphical models,” *Journal of Climate*, vol. 25, no. 17, pp. 5648–5665, 2012.
- [58] EBERT-UPHOFF, I. and DENG, Y., “Causal discovery from spatio-temporal data with applications to climate science,” in *Machine Learning and Applications (ICMLA), 2014 13th International Conference on*, pp. 606–613, IEEE, 2014.
- [59] EVANS, T. and LAMBIOTTE, R., “Line graphs, link partitions, and overlapping communities,” *Physical Review E*, vol. 80, no. 1, p. 016105, 2009.
- [60] FAGHMOUS, J. H. and KUMAR, V., “Spatio-temporal data mining for climate data: Advances, challenges, and opportunities,” in *Data Mining and Knowledge Discovery for Big Data*, pp. 83–116, Springer, 2014.
- [61] FELDHOFF, J. H., LANGE, S., VOLKHOLZ, J., DONGES, J. F., KURTHS, J., and GERSTENGARBE, F.-W., “Complex networks for climate model evaluation with application to statistical versus dynamical modeling of south american climate,” *Climate Dynamics*, vol. 44, no. 5-6, pp. 1567–1581, 2015.

- [62] FOREST, C. E., STONE, P. H., SOKOLOV, A. P., ALLEN, M. R., and WEBSTER, M. D., “Quantifying uncertainties in climate system properties with the use of recent climate observations,” *Science*, vol. 295, no. 5552, pp. 113–117, 2002.
- [63] FORNITO, A., ZALESKY, A., and BREAKSPEAR, M., “The connectomics of brain disorders,” *Nature Reviews Neuroscience*, vol. 16, no. 3, pp. 159–172, 2015.
- [64] FORTUNATO, S., “Community detection in graphs,” *Physics Reports*, vol. 486, no. 3, pp. 75–174, 2010.
- [65] FOUNTALIS, I., BRACCO, A., DILKINA, B., DOVROLIS, C., and KEILHOLZ, S., “ $\{\backslash \text{delta}\}$ -maps: From spatio-temporal data to a weighted and lagged network between functional domains,” *arXiv preprint arXiv:1602.07249*, 2016.
- [66] FOUNTALIS, I., BRACCO, A., and DOVROLIS, C., “Spatio-temporal network analysis for studying climate patterns,” *Climate dynamics*, vol. 42, no. 3-4, pp. 879–899, 2014.
- [67] FOUNTALIS, I., BRACCO, A., and DOVROLIS, C., “Enso in cmip5 simulations: network connectivity from the recent past to the twenty-third century,” *Climate Dynamics*, vol. 45, no. 1-2, pp. 511–538, 2015.
- [68] FOX, M. D., SNYDER, A. Z., VINCENT, J. L., CORBETTA, M., VAN ESSEN, D. C., and RAICHLE, M. E., “The human brain is intrinsically organized into dynamic, anticorrelated functional networks,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 102, no. 27, pp. 9673–9678, 2005.
- [69] FU, K.-S. and MUI, J., “A survey on image segmentation,” *Pattern recognition*, vol. 13, no. 1, pp. 3–16, 1981.
- [70] FYFE, J. C., GILLET, N. P., and ZWIERS, F. W., “Overestimated global warming over the past 20 years,” *Nature Climate Change*, vol. 3, no. 9, pp. 767–769, 2013.
- [71] GAO, J., BULDYREV, S. V., STANLEY, H. E., and HAVLIN, S., “Networks formed from interdependent networks,” *Nature physics*, vol. 8, no. 1, pp. 40–48, 2012.
- [72] GHIL, M. and VAUTARD, R., “Interdecadal oscillations and the warming trend in global temperature time series,” *Nature*, vol. 350, pp. 324–327, 1991.
- [73] GHIL, M., ALLEN, M., DETTINGER, M., IDE, K., KONDRASHOV, D., MANN, M., ROBERTSON, A. W., SAUNDERS, A., TIAN, Y., VARADI, F., and OTHERS, “Advanced spectral methods for climatic time series,” *Reviews of geophysics*, vol. 40, no. 1, 2002.
- [74] GLASSER, M. F., SOTIROPOULOS, S. N., WILSON, J. A., COALSON, T. S., FISCHL, B., ANDERSSON, J. L., XU, J., JBABDI, S., WEBSTER, M., POLIMENI, J. R., and OTHERS, “The minimal preprocessing pipelines for the Human Connectome Project,” *Neuroimage*, vol. 80, pp. 105–124, 2013.

- [75] GORDON, C., COOPER, C., SENIOR, C. A., BANKS, H., GREGORY, J. M., JOHNS, T. C., MITCHELL, J. F., and WOOD, R. A., “The simulation of sst, sea ice extents and ocean heat transports in a version of the hadley centre coupled model without flux adjustments,” *Climate Dynamics*, vol. 16, no. 2-3, pp. 147–168, 2000.
- [76] GORDON, E. M., LAUMANN, T. O., ADEYEMO, B., HUCKINS, J. F., KELLEY, W. M., and PETERSEN, S. E., “Generation and evaluation of a cortical area parcellation from resting-state correlations,” *Cerebral Cortex*, p. bhu239, 2014.
- [77] GOZOLCHIANI, A., HAVLIN, S., and YAMASAKI, K., “Emergence of el niño as an autonomous component in the climate network,” *Physical review letters*, vol. 107, no. 14, p. 148501, 2011.
- [78] GRAHAM, N., “Decadal-scale climate variability in the tropical and north pacific during the 1970s and 1980s: Observations and model results,” *Climate Dynamics*, vol. 10, no. 3, pp. 135–162, 1994.
- [79] GUO, D., “Regionalization with dynamically constrained agglomerative clustering and partitioning (redcap),” *International Journal of Geographical Information Science*, vol. 22, no. 7, pp. 801–823, 2008.
- [80] HANSEN, J., SATO, M., NAZARENKO, L., RUEDY, R., LACIS, A., KOCH, D., TEGEN, I., HALL, T., SHINDELL, D., SANTER, B., and OTHERS, “Climate forcings in goddard institute for space studies si2000 simulations,” *Journal of Geophysical Research: Atmospheres*, vol. 107, no. D18, 2002.
- [81] HELLER, R., STANLEY, D., YEKUTIELI, D., RUBIN, N., and BENJAMINI, Y., “Cluster-based analysis of fmri data,” *NeuroImage*, vol. 33, no. 2, pp. 599–608, 2006.
- [82] HIKOSAKA, K., IWAI, E., SAITO, H., and TANAKA, K., “Polysensory properties of neurons in the anterior bank of the caudal superior temporal sulcus of the macaque monkey,” *Journal of neurophysiology*, vol. 60, no. 5, pp. 1615–1637, 1988.
- [83] HINNE, M., EKMAN, M., JANSSEN, R. J., HESKES, T., and VAN GERVEN, M. A., “Probabilistic clustering of the human connectome identifies communities and hubs,” *PloS one*, vol. 10, no. 1, p. e0117179, 2015.
- [84] HLINKA, J., HARTMAN, D., VEJMEKKA, M., RUNGE, J., MARWAN, N., KURTHS, J., and PALUŠ, M., “Reliability of inference of directed climate networks using conditional mutual information,” *Entropy*, vol. 15, no. 6, pp. 2023–2045, 2013.
- [85] HOLTON, J. R., DMOWSKA, R., and PHILANDER, S. G., *El Niño, La Niña, and the southern oscillation*, vol. 46. Academic press, 1989.
- [86] HUBERT, L. and ARABIE, P., “Comparing partitions,” *Journal of classification*, vol. 2, no. 1, pp. 193–218, 1985.

- [87] HURRELL, J. W. and TRENBERTH, K. E., “Global sea surface temperature analyses: multiple problems and their implications for climate analysis, modeling, and reanalysis,” *Bulletin of the American Meteorological Society*, vol. 80, no. 12, pp. 2661–2678, 1999.
- [88] HYVÄRINEN, A., “Fast and robust fixed-point algorithms for independent component analysis,” *Neural Networks, IEEE Transactions on*, vol. 10, no. 3, pp. 626–634, 1999.
- [89] HYVÄRINEN, A., KARHUNEN, J., and OJA, E., *Independent component analysis*, vol. 46. John Wiley & Sons, 2004.
- [90] JAIN, A. K., MURTY, M. N., and FLYNN, P. J., “Data clustering: a review,” *ACM computing surveys (CSUR)*, vol. 31, no. 3, pp. 264–323, 1999.
- [91] JAIN, N. and DILKINA, B., “Coarse models for bird migrations using clustering and non-stationary markov chains,” in *Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [92] JOLLIFFE, I., *Principal component analysis*. Wiley Online Library, 2002.
- [93] KALNAY, E., KANAMITSU, M., KISTLER, R., COLLINS, W., DEAVEN, D., GANDIN, L., IREDELL, M., SAHA, S., WHITE, G., WOOLLEN, J., and OTHERS, “The ncep/ncar 40-year reanalysis project,” *Bulletin of the American meteorological Society*, vol. 77, no. 3, pp. 437–471, 1996.
- [94] KAWALE, J., CHATTERJEE, S., ORMSBY, D., STEINHAEUSER, K., LIESS, S., and KUMAR, V., “Testing the significance of spatio-temporal teleconnection patterns,” in *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 642–650, ACM, 2012.
- [95] KAWALE, J., LIESS, S., KUMAR, A., STEINBACH, M., GANGULY, A. R., SAMATOVA, N. F., SEMAZZI, F. H., SNYDER, P. K., and KUMAR, V., “Data guided discovery of dynamic climate dipoles,” in *CIDU*, pp. 30–44, 2011.
- [96] KEIL, J. M. and BRECHT, T. B., “The complexity of clustering in planar graphs,” *J. Combinatorial Mathematics and Combinatorial Computing*, vol. 9, pp. 155–159, 1991.
- [97] KIRKMAN IV, C. H. and BITZ, C. M., “The effect of the sea ice freshwater flux on southern ocean temperatures in ccsm3: Deep-ocean warming and delayed surface warming,” *Journal of Climate*, vol. 24, no. 9, pp. 2224–2237, 2011.
- [98] KLEIN, S. A., SODEN, B. J., and LAU, N.-C., “Remote sea surface temperature variations during enso: Evidence for a tropical atmospheric bridge,” *Journal of Climate*, vol. 12, no. 4, pp. 917–932, 1999.
- [99] KOSAKA, Y. and XIE, S.-P., “Recent global-warming hiatus tied to equatorial pacific surface cooling,” *Nature*, vol. 501, no. 7467, pp. 403–407, 2013.

- [100] KRAMER, M. A., EDEN, U. T., CASH, S. S., and KOLACZYK, E. D., “Network inference with confidence from multivariate time series,” *Physical Review E*, vol. 79, no. 6, p. 061916, 2009.
- [101] KRETSCHMER, M., COUMOU, D., DONGES, J. F., and RUNGE, J., “Using causal effect networks to analyze different arctic drivers of mid-latitude winter circulation,” *Journal of Climate*, no. 2016, 2016.
- [102] KUCHARSKI, F., KANG, I.-S., FARNETI, R., and FEUDALE, L., “Tropical pacific response to 20th century atlantic warming,” *Geophysical Research Letters*, vol. 38, no. 3, 2011.
- [103] LANCICHINETTI, A., RADICCHI, F., RAMASCO, J. J., FORTUNATO, S., and OTHERS, “Finding statistically significant communities in networks,” *PloS one*, vol. 6, no. 4, p. e18961, 2011.
- [104] LU, Y., JIANG, T., and ZANG, Y., “Region growing method for the analysis of functional mri data,” *NeuroImage*, vol. 20, no. 1, pp. 455–465, 2003.
- [105] LUDESCHER, J., GOZOLCHIANI, A., BOGACHEV, M. I., BUNDE, A., HAVLIN, S., and SCHELLNHUBER, H. J., “Improved el niño forecasting by cooperativity detection,” *Proceedings of the National Academy of Sciences*, vol. 110, no. 29, pp. 11742–11745, 2013.
- [106] MALIK, N., BOOKHAGEN, B., MARWAN, N., and KURTHS, J., “Analysis of spatial and temporal extreme monsoonal rainfall over south asia using complex networks,” *Climate dynamics*, vol. 39, no. 3-4, pp. 971–987, 2012.
- [107] MARTIN, E. and DAVIDSEN, J., “Estimating time delays for constructing dynamical networks,” *Nonlinear Processes in Geophysics*, vol. 21, no. 5, pp. 929–937, 2014.
- [108] MCGUIRE, M. P. and NGUYEN, N. P., “Community structure analysis in big climate data,” in *Big Data (Big Data), 2014 IEEE International Conference on*, pp. 38–46, IEEE, 2014.
- [109] MEINSHAUSEN, M., SMITH, S. J., CALVIN, K., DANIEL, J. S., KAINUMA, M., LAMARQUE, J., MATSUMOTO, K., MONTZKA, S., RAPER, S., RIAHI, K., and OTHERS, “The rcp greenhouse gas concentrations and their extensions from 1765 to 2300,” *Climatic change*, vol. 109, no. 1-2, pp. 213–241, 2011.
- [110] MILLER, A. J., CAYAN, D. R., BARNETT, T. P., GRAHAM, N. E., and OBERHUBER, J. M., “The 1976–77 climate shift of the pacific ocean,” *Oceanography*, vol. 7, no. 1, pp. 21–26, 1994.
- [111] MORENO-DOMINGUEZ, D., ANWANDER, A., and KNÖSCHE, T. R., “A hierarchical method for whole-brain connectivity-based parcellation,” *Human brain mapping*, vol. 35, no. 10, pp. 5000–5025, 2014.

- [112] NADLER, B. and GALUN, M., “Fundamental limitations of spectral clustering,” in *Advances in Neural Information Processing Systems*, pp. 1017–1024, 2006.
- [113] NEWMAN, M., *Networks: an introduction*. OUP Oxford, 2010.
- [114] NEWMAN, M., BARABASI, A.-L., and WATTS, D. J., *The structure and dynamics of networks*. Princeton University Press, 2006.
- [115] NEWMAN, M. E. and GIRVAN, M., “Finding and evaluating community structure in networks,” *Physical review E*, vol. 69, no. 2, p. 026113, 2004.
- [116] PALLA, G., DERÉNYI, I., FARKAS, I., and VICSEK, T., “Uncovering the overlapping community structure of complex networks in nature and society,” *Nature*, vol. 435, no. 7043, pp. 814–818, 2005.
- [117] PELAN, A., STEINHAEUSER, K., CHAWLA, N. V., DE ALWIS PITTS, D. A., and GANGULY, A. R., “Empirical comparison of correlation measures and pruning levels in complex networks representing the global climate system,” in *Computational Intelligence and Data Mining (CIDM), 2011 IEEE Symposium on*, pp. 239–245, IEEE, 2011.
- [118] POWER, J. D., COHEN, A. L., NELSON, S. M., WIG, G. S., BARNES, K. A., CHURCH, J. A., VOGEL, A. C., LAUMANN, T. O., MIEZIN, F. M., SCHLAGGAR, B. L., and OTHERS, “Functional network organization of the human brain,” *Neuron*, vol. 72, no. 4, pp. 665–678, 2011.
- [119] RADEBACH, A., DONNER, R. V., RUNGE, J., DONGES, J. F., and KURTHS, J., “Disentangling different types of el niño episodes by evolving climate network analysis,” *Physical Review E*, vol. 88, no. 5, p. 052807, 2013.
- [120] RAND, W. M., “Objective criteria for the evaluation of clustering methods,” *Journal of the American Statistical association*, vol. 66, no. 336, pp. 846–850, 1971.
- [121] RAYNER, N., PARKER, D. E., HORTON, E., FOLLAND, C., ALEXANDER, L., ROWELL, D., KENT, E., and KAPLAN, A., “Global analyses of sea surface temperature, sea ice, and night marine air temperature since the late nineteenth century,” *Journal of Geophysical Research: Atmospheres (1984–2012)*, vol. 108, no. D14, 2003.
- [122] RESHEF, D. N., RESHEF, Y. A., FINUCANE, H. K., GROSSMAN, S. R., MCVEAN, G., TURNBAUGH, P. J., LANDER, E. S., MITZENMACHER, M., and SABETI, P. C., “Detecting novel associations in large data sets,” *science*, vol. 334, no. 6062, pp. 1518–1524, 2011.
- [123] REYNOLDS, R. W. and SMITH, T. M., “Improved global sea surface temperature analyses using optimum interpolation,” *Journal of climate*, vol. 7, no. 6, pp. 929–948, 1994.

- [124] RIAHI, K., RAO, S., KREY, V., CHO, C., CHIRKOV, V., FISCHER, G., KINDERMANN, G., NAKICENOVIC, N., and RAFAJ, P., “Rcp 8.5a scenario of comparatively high greenhouse gas emissions,” *Climatic Change*, vol. 109, no. 1-2, pp. 33–57, 2011.
- [125] RODRÍGUEZ-FONSECA, B., POLO, I., GARCÍA-SERRANO, J., LOSADA, T., MOHINO, E., MECHOSO, C. R., and KUCHARSKI, F., “Are atlantic niños enhancing pacific enso events in recent decades?,” *Geophysical Research Letters*, vol. 36, no. 20, 2009.
- [126] ROGERS, G. S., “A course in theoretical statistics,” *Technometrics*, vol. 11, no. 4, pp. 840–841, 1969.
- [127] RUMMEL, C., MÜLLER, M., BAIER, G., AMOR, F., and SCHINDLER, K., “Analyzing spatio-temporal patterns of genuine cross-correlations,” *Journal of neuroscience methods*, vol. 191, no. 1, pp. 94–100, 2010.
- [128] RUNGE, J., PETOUKHOV, V., DONGES, J. F., HLINKA, J., JAJCAY, N., VEMELKA, M., HARTMAN, D., MARWAN, N., PALUŠ, M., and KURTHS, J., “Identifying causal gateways and mediators in complex spatio-temporal systems,” *Nature communications*, vol. 6, 2015.
- [129] SANTOSO, A., MCGREGOR, S., JIN, F.-F., CAI, W., ENGLAND, M. H., AN, S.-I., MCPHADEN, M. J., and GUILYARDI, E., “Late-twentieth-century emergence of the el niño propagation asymmetry and future projections,” *Nature*, vol. 504, no. 7478, pp. 126–130, 2013.
- [130] SCHLÖSSER, R., GESIERICH, T., KAUFMANN, B., VUCUREVIC, G., HUNSCH, S., GAWEHN, J., and STOETER, P., “Altered effective connectivity during working memory performance in schizophrenia: a study with fmri and structural equation modeling,” *Neuroimage*, vol. 19, no. 3, pp. 751–763, 2003.
- [131] SHEN, X., TOKOGLU, F., PAPADEMETRIS, X., and CONSTABLE, R. T., “Group-wise whole-brain parcellation from resting-state fmri data for network node identification,” *Neuroimage*, vol. 82, pp. 403–415, 2013.
- [132] SIMMONS, A., WALLACE, J., and BRANSTATOR, G., “Barotropic wave propagation and instability, and atmospheric teleconnection patterns,” *Journal of the Atmospheric Sciences*, vol. 40, no. 6, pp. 1363–1392, 1983.
- [133] SMITH, S. M., BECKMANN, C. F., ANDERSSON, J., AUERBACH, E. J., BIJSTERBOSCH, J., DOUAUD, G., DUFF, E., FEINBERG, D. A., GRIFFANTI, L., HARMS, M. P., and OTHERS, “Resting-state fMRI in the human connectome project,” *Neuroimage*, vol. 80, pp. 144–168, 2013.
- [134] SMITH, T. M., REYNOLDS, R. W., PETERSON, T. C., and LAWRIK, J., “Improvements to noaa’s historical merged land-ocean surface temperature analysis (1880-2006),” *Journal of Climate*, vol. 21, no. 10, pp. 2283–2296, 2008.

- [135] SOLOMON, A. and NEWMAN, M., “Reconciling disparate twentieth-century indopacific ocean temperature trends in the instrumental record,” *Nature Climate Change*, vol. 2, no. 9, pp. 691–699, 2012.
- [136] SPORNS, O., *Networks of the Brain*. MIT press, 2011.
- [137] SPORNS, O. and BETZEL, R. F., “Modular brain networks,” *Annual review of psychology*, vol. 67, no. 1, 2015.
- [138] STAM, C. J., “Modern network science of neurological disorders,” *Nature Reviews Neuroscience*, vol. 15, no. 10, pp. 683–695, 2014.
- [139] STEINBACH, M., TAN, P.-N., KUMAR, V., KLOOSTER, S., and POTTER, C., “Discovery of climate indices using clustering,” in *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 446–455, ACM, 2003.
- [140] STEINHAEUSER, K. and CHAWLA, N. V., “Identifying and evaluating community structure in complex networks,” *Pattern Recognition Letters*, vol. 31, no. 5, pp. 413–421, 2010.
- [141] STEINHAEUSER, K., CHAWLA, N. V., and GANGULY, A. R., “Complex networks in climate science: Progress, opportunities and challenges,” in *CIDU*, pp. 16–26, 2010.
- [142] STEINHAEUSER, K., CHAWLA, N. V., and GANGULY, A. R., “An exploration of climate data using complex networks,” *ACM SIGKDD Explorations Newsletter*, vol. 12, no. 1, pp. 25–32, 2010.
- [143] STEINHAEUSER, K., CHAWLA, N. V., and GANGULY, A. R., “Complex networks as a unified framework for descriptive analysis and predictive modeling in climate science,” *Statistical Analysis and Data Mining*, vol. 4, no. 5, pp. 497–511, 2011.
- [144] STEINHAEUSER, K., GANGULY, A. R., and CHAWLA, N. V., “Multivariate and multiscale dependence in the global climate system revealed through complex networks,” *Climate dynamics*, vol. 39, no. 3-4, pp. 889–895, 2012.
- [145] STEINHAEUSER, K. and TSONIS, A. A., “A climate model intercomparison at the dynamics level,” *Climate dynamics*, vol. 42, no. 5-6, pp. 1665–1670, 2014.
- [146] STEVENS, B., BONY, S., and OTHERS, “What are climate models missing,” *Science*, vol. 340, no. 6136, pp. 1053–1054, 2013.
- [147] SUPEKAR, K., MENON, V., RUBIN, D., MUSEN, M., and GREICIUS, M. D., “Network analysis of intrinsic functional brain connectivity in alzheimer’s disease,” *PLoS Comput Biol*, vol. 4, no. 6, p. e1000100, 2008.
- [148] SWANSON, K. L. and TSONIS, A. A., “Has the climate recently shifted?,” *Geophysical Research Letters*, vol. 36, no. 6, 2009.

- [149] TAYLOR, K. E., STOUFFER, R. J., and MEEHL, G. A., “An overview of cmip5 and the experiment design,” *Bulletin of the American Meteorological Society*, vol. 93, no. 4, pp. 485–498, 2012.
- [150] THIRION, B., VAROQUAUX, G., DOHMATOB, E., and POLINE, J.-B., “Which fmri clustering gives good brain parcellations?,” *Frontiers in neuroscience*, vol. 8, 2014.
- [151] THOMPSON, G. J., MERRITT, M. D., PAN, W.-J., MAGNUSON, M. E., GROOMS, J. K., JAEGER, D., and KEILHOLZ, S. D., “Neural correlates of time-varying functional connectivity in the rat,” *Neuroimage*, vol. 83, pp. 826–836, 2013.
- [152] TONONI, G., MCINTOSH, A. R., RUSSELL, D. P., and EDELMAN, G. M., “Functional clustering: identifying strongly interactive brain regions in neuroimaging data,” *Neuroimage*, vol. 7, no. 2, pp. 133–149, 1998.
- [153] TSONIS, A. A. and ROEBBER, P. J., “The architecture of the climate network,” *Physica A: Statistical Mechanics and its Applications*, vol. 333, pp. 497–504, 2004.
- [154] TSONIS, A. A., SWANSON, K., and KRAVTSOV, S., “A new dynamical mechanism for major climate shifts,” *Geophysical Research Letters*, vol. 34, no. 13, 2007.
- [155] TSONIS, A. A. and SWANSON, K. L., “Topology and predictability of el nino and la nina networks,” *Physical Review Letters*, vol. 100, no. 22, p. 228502, 2008.
- [156] TSONIS, A. A., SWANSON, K. L., and ROEBBER, P. J., “What do networks have to do with climate?,” *Bulletin of the American Meteorological Society*, vol. 87, no. 5, p. 585, 2006.
- [157] TSONIS, A. A., SWANSON, K. L., and WANG, G., “On the role of atmospheric teleconnections in climate,” *Journal of Climate*, vol. 21, no. 12, pp. 2990–3001, 2008.
- [158] TSONIS, A. A., WANG, G., SWANSON, K. L., RODRIGUES, F. A., and DA FONTURA COSTA, L., “Community structure and dynamics in climate networks,” *Climate dynamics*, vol. 37, no. 5-6, pp. 933–940, 2011.
- [159] UPPALA, S. M., KÅLLBERG, P., SIMMONS, A., ANDRAE, U., BECHTOLD, V. D., FIORINO, M., GIBSON, J., HASELER, J., HERNANDEZ, A., KELLY, G., and OTHERS, “The era-40 re-analysis,” *Quarterly Journal of the Royal Meteorological Society*, vol. 131, no. 612, pp. 2961–3012, 2005.
- [160] VAN DEN HEUVEL, M., MANDL, R., and POL, H. H., “Normalized cut group clustering of resting-state fmri data,” *PloS one*, vol. 3, no. 4, p. e2001, 2008.
- [161] VAN DEN HEUVEL, M. P. and SPORNS, O., “Rich-club organization of the human connectome,” *The Journal of neuroscience*, vol. 31, no. 44, pp. 15775–15786, 2011.

- [162] VAN ESSEN, D. C., GLASSER, M. F., DIERKER, D. L., HARWELL, J., and COALSON, T., “Parcellations and hemispheric asymmetries of human cerebral cortex analyzed on surface-based atlases,” *Cerebral Cortex*, vol. 22, no. 10, pp. 2241–2262, 2012.
- [163] VAN ESSEN, D. C., SMITH, S. M., BARCH, D. M., BEHRENS, T. E., YACOUB, E., UGURBIL, K., CONSORTIUM, W.-M. H., and OTHERS, “The wu-minn human connectome project: an overview,” *Neuroimage*, vol. 80, pp. 62–79, 2013.
- [164] VEJMEĽKA, M., POKORNÁ, L., HLINKA, J., HARTMAN, D., JAJCAY, N., and PALUŠ, M., “Non-random correlation structures and dimensionality reduction in multivariate climate data,” *Climate Dynamics*, vol. 44, no. 9-10, pp. 2663–2682, 2014.
- [165] VÉRTES, P. E., ALEXANDER-BLOCH, A. F., GOGTAY, N., GIEDD, J. N., RAPOPORT, J. L., and BULLMORE, E. T., “Simple models of human brain functional networks,” *Proceedings of the National Academy of Sciences*, vol. 109, no. 15, pp. 5868–5873, 2012.
- [166] VIDARD, A., ANDERSON, D. L., and BALMASEDA, M., “Impact of ocean observation systems on ocean analysis and seasonal forecasts,” *Monthly weather review*, vol. 135, no. 2, pp. 409–429, 2007.
- [167] VON STORCH, H. and ZWIERS, F. W., *Statistical analysis in climate research*. Cambridge university press, 2001.
- [168] WANG, G., SWANSON, K. L., and TSONIS, A. A., “The pacemaker of major climate shifts,” *Geophysical Research Letters*, vol. 36, no. 7, 2009.
- [169] WARD JR, J. H., “Hierarchical grouping to optimize an objective function,” *Journal of the American statistical association*, vol. 58, no. 301, pp. 236–244, 1963.
- [170] WIG, G. S., LAUMANN, T. O., and PETERSEN, S. E., “An approach for parcellating human cortical areas using resting-state correlations,” *Neuroimage*, vol. 93, pp. 276–291, 2014.
- [171] WU, K., TAKI, Y., SATO, K., SASSA, Y., INOUE, K., GOTO, R., OKADA, K., KAWASHIMA, R., HE, Y., EVANS, A. C., and OTHERS, “The overlapping community structure of structural brain network in young healthy individuals,” *PLoS One*, vol. 6, no. 5, p. e19608, 2011.
- [172] XIE, P. and ARKIN, P. A., “Global precipitation: A 17-year monthly analysis based on gauge observations, satellite estimates, and numerical model outputs,” *Bulletin of the American Meteorological Society*, vol. 78, no. 11, pp. 2539–2558, 1997.
- [173] YAMASAKI, K., GOZOLCHIANI, A., and HAVLIN, S., “Climate networks around the globe are significantly affected by el nino,” *Physical review letters*, vol. 100, no. 22, p. 228501, 2008.

-
- [174] YAMASAKI, K., GOZOLCHIANI, A., and HAVLIN, S., “Climate networks based on phase synchronization analysis track el-nino,” *Progress of Theoretical Physics Supplement*, vol. 179, pp. 178–188, 2009.
- [175] YAN, X., KELLEY, S., GOLDBERG, M., and BISWAL, B. B., “Detecting overlapped functional clusters in resting state fmri with connected iterative scan: a graph theory based clustering algorithm,” *Journal of neuroscience methods*, vol. 199, no. 1, pp. 108–118, 2011.
- [176] YEO, B. T., KRIENEN, F. M., SEPULCRE, J., SABUNCU, M. R., LASHKARI, D., HOLLINSHEAD, M., ROFFMAN, J. L., SMOLLER, J. W., ZÖLLEI, L., POLIMENI, J. R., and OTHERS, “The organization of the human cerebral cortex estimated by intrinsic functional connectivity,” *Journal of neurophysiology*, vol. 106, no. 3, pp. 1125–1165, 2011.
- [177] ZHANG, P., HUANG, Y., SHEKHAR, S., and KUMAR, V., “Correlation analysis of spatial time series datasets: A filter-and-refine approach,” in *Advances in Knowledge Discovery and Data Mining*, pp. 532–544, Springer, 2003.
- [178] ZHANG, W. and JIN, F.-F., “Improvements in the cmip5 simulations of enso-ssta meridional width,” *Geophysical Research Letters*, vol. 39, no. 23, 2012.
- [179] ZHANG, W., JIN, F.-F., ZHAO, J.-X., and LI, J., “On the bias in simulated enso ssta meridional widths of cmip3 models,” *Journal of Climate*, vol. 26, no. 10, pp. 3173–3186, 2013.