# Simplify to understand: extracting insight from high-dimensional simulation datasets of biological systems via information-theoretic approaches

#### **Roberto Menichetti**



Physics Department, University of Trento, Italy Trento Institute for Fundamental Physics and Applications (INFN-TIFPA), Italy

> ALPACA Workshop ECT\*, Trento, September 2023



Statistical & Biological Physics





#### Data are easier and easier to get



1







#### All-atom MD simulations



Hadden *et al.*, eLife 2018



#### Data are easier and easier to get







#### **All-atom MD simulations**



Hadden et al., eLife 2018



#### Data are easier and easier to get

#### Aggregate digital databases



Morand, Yip, Velegrakis, Lattanzi, Potestio, Tubiana, arXiv 2023







#### **All-atom MD simulations**



Hadden et al., eLife 2018



#### Data are easier and easier to get

#### Aggregate digital databases

#### Huge amount of generated data: how to extract insight?

Morand, Yip, Velegrakis, Lattanzi, Potestio, Tubiana, arXiv 2023





### Data are easier and easier to get

#### All-atom MD simulations



Hadden *et al.*, eLife 2018



 $\vec{F} = m\vec{a}$ 10<sup>6</sup> atoms for 10<sup>9</sup> configurations of the system

ATOM	I 1153	CD	GLN	Х	408	110.020	92.423	42.297	0.00	0.00
ATOM	1 1154	OE1	GLN	Х	408	110.283	93.074	43.331	0.00	0.00
ATOM	I 1155	NE2	GLN	Х	408	108.983	91.618	42.373	0.00	0.00
ATOM	I 1156	HE21	GLN	Х	408	108.656	91.204	41.511	0.00	0.00
ATOM	1 1157	HE22	GLN	Х	408	108.471	91.768	43.225	0.00	0.00
ATOM	1 1158	С	GLN	Х	408	113.167	95.332	40.798	0.00	0.00
ATOM	I 1159	0	GLN	Х	408	113.764	95.563	39.715	0.00	0.00
ATOM	1160 I	N	GLN	Х	409	113.578	95.654	42.042	0.00	0.00
ATOM	1161 I	Н	GLN	Х	409	113.018	95.252	42.780	0.00	0.00
ATOM	I 1162	CA	GLN	Х	409	114.881	96.418	42.260	0.00	0.00
ATOM	I 1163	HA	GLN	Х	409	115.146	96.965	41.356	0.00	0.00
ATOM	1164	CB	GLN	Х	409	114.444	97.423	43.339	0.00	0.00
ATOM	I 1165	HB1	GLN	Х	409	115.314	97.840	43.844	0.00	0.00
ATOM	1 1166	HB2	GLN	Х	409	113.900	96.872	44.111	0.00	0.00
ATOM	1 1167	CG	GLN	Х	409	113.759	98.671	42.779	0.00	0.00
ATOM	I 1168	HG1	GLN	Х	409	112.955	98.388	42.101	0.00	0.00
ATOM	1169	HG2	GLN	Х	409	114.393	99.253	42.117	0.00	0.00
ATOM	I 1170	CD	GLN	Х	409	113.157	99.520	43.821	0.00	0.00
ATOM	I 1171	OE1	GLN	Х	409	113.911	100.138	44.580	0.00	0.00
ATOM	I 1172	NE2	GLN	Х	409	111.814	99.753	43.795	0.00	0.00
ATOM	I 1173	HE21	GLN	Х	409	111.336	99.317	43.023	0.00	0.00
ATOM	1 1174	HE22	GLN	Х	409	111.343	100.484	44.321	0.00	0.00
ATOM	1 1175	С	GLN	Х	409	116.132	95.610	42.722	0.00	0.00
ATOM	1 1176	0	GLN	Х	409	117.199	96.174	42.994	0.00	0.00
ATOM	I 1177	Ν	PHE	Х	410	115.988	94.277	42.829	0.00	0.00
ATOM	I 1178	Н	PHE	Х	410	115.077	93.891	42.653	0.00	0.00
ATOM	1 1179	CA	PHE	Х	410	116.976	93.435	43.590	0.00	0.00
ATOM	1180	HA	PHE	Х	410	117.124	93.938	44.557	0.00	0.00
ATOM	I 1181	CB	PHE	Х	410	116.299	92.042	43.824	0.00	0.00
ATOM	I 1182	HB1	PHE	Х	410	116.161	91.560	42.866	0.00	0.00
ATOM	I 1183	HB2	PHE	Х	410	115.259	92.135	44.154	0.00	0.00
ATOM	1 1184	CG	PHE	Х	410	117.011	91.074	44.747	0.00	0.00
ATOM	I 1185	CD1	PHE	Х	410	117.099	91.369	46.110	0.00	0.00
ATOM	1 1186	HD1	PHE	Х	410	116.520	92.142	46.584	0.00	0.00
ATOM	I 1187	CE1	PHE	Х	410	117.746	90.542	46.998	0.00	0.00
ATOM	I 1188	HE1	PHE	Х	410	117.720	90.932	48.007	0.00	0.00
ATOM	I 1189	СZ	PHE	Х	410	118.393	89.398	46.545	0.00	0.00
ATOM	1 1190	ΗZ	PHE	Х	410	118.896	88.709	47.208	0.00	0.00
ATOM	1 1191	CE2	PHE	Х	410	118.374	89.090	45.183	0.00	0.00
ATOM	1 1192	HE2	PHE	Х	410	118.911	88.260	44.740	0.00	0.00
ATOM	1 1193	CD2	PHE	Х	410	117.779	89.966	44.293	0.00	0.00
ATOM	1194	HD2	PHE	Х	410	117.894	89.832	43.227	0.00	0.00
ATOM	1 1195	С	PHE	Х	410	118.379	93.352	42.972	0.00	0.00















#### DATA != INFORMATION













#### DATA != INFORMATION



Simplify "just right" to understand without losing information











#### DATA != INFORMATION



Simplify "just right" to understand without losing information

#### Simplify to understand





Dimensionality reduction: in terms of what variables?





ATOM	1153	CD	GLN	Х	408	110.020	92.423	42.297	0.00	0.00
ATOM	1154	OE1	GLN	Х	408	110.283	93.074	43.331	0.00	0.00
ATOM	1155	NE2	GLN	Х	408	108.983	91.618	42.373	0.00	0.00
ATOM	1156	HE21	GLN	Х	408	108.656	91.204	41.511	0.00	0.00
ATOM	1157	HE22	GLN	Х	408	108.471	91.768	43.225	0.00	0.00
ATOM	1158	С	GLN	Х	408	113.167	95.332	40.798	0.00	0.00
ATOM	1159	0	GLN	Х	408	113.764	95.563	39.715	0.00	0.00
ATOM	1160	Ν	GLN	Х	409	113.578	95.654	42.042	0.00	0.00
ATOM	1161	Η	GLN	Х	409	113.018	95.252	42.780	0.00	0.00
ATOM	1162	CA	GLN	Х	409	114.881	96.418	42.260	0.00	0.00
ATOM	1163	HA	GLN	Х	409	115.146	96.965	41.356	0.00	0.00
ATOM	1164	CB	GLN	Х	409	114.444	97.423	43.339	0.00	0.00
ATOM	1165	HB1	GLN	Х	409	115.314	97.840	43.844	0.00	0.00
ATOM	1166	HB2	GLN	Х	409	113.900	96.872	44.111	0.00	0.00









# Simplify to understand

ATOM ATOM ATOM ATOM ATOM ATOM ATOM	1153 CD 1154 OE1 1155 NE2 1156 HE21 1157 HE22 1158 C 1159 O 1160 N	GLN X 408 GLN X 408	110.020 110.283 108.983 108.656 108.471 113.167 113.764	92.423 93.074 91.618 91.204 91.768 95.332 95.563	42.297 43.331 42.373 41.511 43.225 40.798 39.715	0.00 0.00 0.00 0.00 0.00 0.00 0.00	0.00 0.00 0.00 0.00 0.00 0.00 0.00	
ATOM ATOM ATOM ATOM ATOM ATOM	1161 H 1162 CA 1163 HA 1164 CB 1165 HB1 1166 HB2	GLN X 409 GLN X 409 GLN X 409 GLN X 409 GLN X 409 GLN X 409 GLN X 409	113.018 114.881 115.146 114.444 115.314 113.900	95 96 96 97 97 97	200 -			
				) 	100 - <b>-</b> 50 -			
		R			0 - 50 -			
					-200	-150	-100	-50 <b>Φ</b>



**Collective variables** 







ATOM ATOM ATOM ATOM ATOM ATOM ATOM	1153 CD 1154 OE1 1155 NE2 1156 HE21 1157 HE22 1158 C 1159 O 1160 N	GLN X 408 GLN X 408	110.020 110.283 108.983 108.656 108.471 113.167 113.764	92.423 93.074 91.618 91.204 91.768 95.332 95.563	42.297 43.331 42.373 41.511 43.225 40.798 39.715	0.00 0.00 0.00 0.00 0.00 0.00 0.00	0.00 0.00 0.00 0.00 0.00 0.00 0.00	
ATOM ATOM ATOM ATOM ATOM ATOM	1161 H 1162 CA 1163 HA 1164 CB 1165 HB1 1166 HB2	GLN X 409 GLN X 409 GLN X 409 GLN X 409 GLN X 409 GLN X 409 GLN X 409	113.018 114.881 115.146 114.444 115.314 113.900	95 96 96 97 97 97	200 -			
				) 	100 - <b>-</b> 50 -			
		R			0 - 50 -			
					-200	-150	-100	-50 <b>Φ</b>







#### AA representation



Rudzinski, Noid, J. Chem. Phys. 2011



# **Coarse-graining (CG)**







#### CG representation



Rudzinski, Noid, J. Chem. Phys. 2011



# **Coarse-graining (CG)**

#### Describe the system in terms of a reduced number of particles: CG filter







#### CG representation



Rudzinski, Noid, J. Chem. Phys. 2011



# **Coarse-graining (CG)**

#### Describe the system in terms of a reduced number of particles: CG filter

Degree of detail employed can be **modulated** 

**Intuitive** interpretation of the outcome



































# High-resolution model





### **CG filters**

Choice of the sites











**Observe** the system with a pair of CG glasses: information loss











**Observe** the system with a pair of CG glasses: information loss

Optimize the glasses to extract relevant information from the data











#### Observe the system with a pair of CG glasses: information loss

**Optimize** the glasses to extract relevant information from the data

# understand











$$\chi^{2} = \frac{1}{3N} \sum_{\mathrm{I}=1}^{\mathrm{N}} \frac{1}{n_{\mathrm{t}}} \sum_{\mathrm{t}=1}^{\mathrm{n}_{\mathrm{t}}} \left( \sum_{\mathrm{i}\in\mathrm{I}} \sum_{\mathrm{j}\geq\mathrm{i}\in\mathrm{I}} |\Delta \mathbf{r}_{\mathrm{i}}^{\mathrm{ED}}(t) - \Delta \mathbf{r}_{\mathrm{j}}^{\mathrm{ED}}(t)|^{2} \right)$$

Zhang, Lu, Noid, Krishna, Pfaendtner, Voth BJ 2008



Potestio, Pontiggia, Micheletti BJ 2009



### **CG** filters



# understand















$$\chi^2 = \frac{1}{3N} \sum_{I=1}^{N} \frac{1}{n_t} \sum_{t=1}^{n_t} \left( \sum_{i \in I} \sum_{j \ge i \in I} |\Delta \mathbf{r}_i^{ED}(t) - \Delta \mathbf{r}_j^{ED}(t) \right)$$

Zhang, Lu, Noid, Krishna, Pfaendtner, Vol BJ 2008

Optimize the CG filter based on quantitative criteria: In our case, minimize the information loss arising from coarse-graining



Potestio, Pontiggia, Micheletti BJ 2009



## **CG filters**

# understand





















Rudzinski, Noid, J. Chem. Phys. 2011 Giulini, RM, Shell, Potestio, JCTC 2020



### Information loss in coarse-graining









 $\bar{p}_r(\mathbf{r}) \neq p_r(\mathbf{r})$ 









![](_page_27_Picture_6.jpeg)

![](_page_28_Picture_1.jpeg)

![](_page_28_Picture_2.jpeg)

![](_page_28_Picture_4.jpeg)

![](_page_28_Picture_7.jpeg)

![](_page_29_Picture_1.jpeg)

![](_page_29_Figure_2.jpeg)

Original high resolution representation

Rudzinski, Noid, J. Chem. Phys. 2011 Giulini, RM, Shell, Potestio, JCTC 2020

![](_page_29_Picture_5.jpeg)

![](_page_29_Figure_6.jpeg)

Each microstate **r** receives the average probability of the macrostate **R** it maps onto

![](_page_29_Figure_8.jpeg)

Filtering

Reconstructed high resolution representation

![](_page_29_Picture_12.jpeg)

![](_page_30_Picture_1.jpeg)

![](_page_30_Figure_2.jpeg)

Original high resolution representation

Filtering

Rudzinski, Noid, J. Chem. Phys. 2011 Giulini, RM, Shell, Potestio, JCTC 2020

![](_page_30_Picture_6.jpeg)

Reconstructed high resolution representation

![](_page_30_Picture_10.jpeg)

![](_page_31_Picture_1.jpeg)

![](_page_31_Figure_2.jpeg)

Rudzinski, Noid, J. Chem. Phys. 2011 Giulini, RM, Shell, Potestio, JCTC 2020

![](_page_31_Picture_4.jpeg)

![](_page_31_Figure_5.jpeg)

 $\bar{p}_r(\mathbf{r}) = \frac{p_R(\mathbf{M}(\mathbf{r}))}{\Sigma(\mathbf{M}(\mathbf{r}))}$  Each microstate **r** receives the average probability of the macrostate **R** it maps onto

Mapping entropy

$$\mathbf{r} \ p_r(\mathbf{r}) \ln \left[ \frac{p_r(\mathbf{r})}{\bar{p}_r(\mathbf{r})} \right]$$

$$\left[\frac{p_r(\mathbf{r})}{\bar{p}_r(\mathbf{r})}\right]$$

Quantifies the information loss arising from coarse-graining **DEPENDS ONLY ON THE CHOICE OF THE CG SITES** 

Filtering

high resolution representation

![](_page_31_Picture_14.jpeg)

![](_page_32_Picture_1.jpeg)

 $S_{map} \approx k_B \frac{\beta^2}{2} \int d\mathbf{R} p_R(\mathbf{R}) \langle (u - \langle u \rangle_{\beta|\mathbf{R}})^2 \rangle_{\beta|\mathbf{R}}$ 

Giulini, RM, Shell, Potestio, JCTC 2020

![](_page_32_Picture_4.jpeg)

![](_page_32_Picture_7.jpeg)

![](_page_33_Picture_1.jpeg)

 $S_{map} \approx k_B \frac{\beta^2}{2} \int d\mathbf{R} p_R(\mathbf{R}) \langle (u - \langle u \rangle_{\beta | \mathbf{R}})^2 \rangle_{\beta | \mathbf{R}}$ 

![](_page_33_Figure_4.jpeg)

Giulini, RM, Shell, Potest

![](_page_33_Picture_6.jpeg)

Calculating the mapping entropy of a CG representation

![](_page_33_Picture_8.jpeg)

![](_page_34_Picture_1.jpeg)

 $S_{map} \approx k_B \frac{\beta^2}{2} \int d\mathbf{R} p_R(\mathbf{R}) \langle (u - \langle u \rangle_{\beta|\mathbf{R}})^2 \rangle_{\beta|\mathbf{R}}$ 

![](_page_34_Picture_4.jpeg)

Giulini, RM, Shell, Potest

![](_page_34_Picture_6.jpeg)

Calculating the mapping entropy of a CG representation

![](_page_34_Picture_8.jpeg)

![](_page_35_Picture_1.jpeg)

 $S_{map} \approx k_B \frac{\beta^2}{2} \int d\mathbf{R}_I$ 

![](_page_35_Figure_4.jpeg)

Giulini, RM, Shell, Potest

![](_page_35_Picture_6.jpeg)

$$\Delta p_R(\mathbf{R})\langle (u - \langle u \rangle_{\beta|\mathbf{R}})^2 \rangle_{\beta|\mathbf{R}}$$

Calculating the mapping entropy of a CG representation

![](_page_35_Picture_10.jpeg)
# Estimating the information loss





Giulini, RM, Shell, Potest







#### Looking for maximally informative selections of CG sites: Stochastic minimization of S<sub>map</sub> over the mapping space















#### Looking for maximally informative selections of CG sites: Stochastic minimization of S<sub>map</sub> over the mapping space









#### Looking for maximally informative selections of CG sites: Stochastic minimization of S<sub>map</sub> over the mapping space









#### Application to Tamapin





Giulini, RM, Shell, Potestio, JCTC 2020







#### Application to Tamapin



**>** 723





#### Application to Tamapin



**<** 723





#### Application to Tamapin







#### Application to Tamapin



#### Application to CzrA transcription repressor









#### Application to Tamapin



#### Application to CzrA transcription repressor







#### Application to Tamapin



#### Application to CzrA transcription repressor







#### Application to Tamapin



#### Application to CzrA transcription repressor











#### Calculating the mapping entropy is a computationally expensive task



Errica, Giulini, RM et al., Front. Mol. Biosci. 2021









#### Calculating the mapping entropy is a computationally expensive task



Errica, Giulini, RM et al., Front. Mol. Biosci. 2021











#### Artificial Intelligence to speed up calculations: Deep Graph Networks

Errica, Giulini, RM et al., Front. Mol. Biosci. 2021









#### Artificial Intelligence to speed up calculations: Deep Graph Networks

Encode the (static) structure of the molecule in a graph



Errica, Giulini, RM et al., Front. Mol. Biosci. 2021







#### Artificial Intelligence to speed up calculations: Deep Graph Networks

#### Encode the (static) structure of the molecule in a graph

Employ DGNs to predict the mapping entropy





Errica, Giulini, RM et al., Front. Mol. Biosci. 2021









#### Artificial Intelligence to speed up calculations: Deep Graph Networks

#### Encode the (static) structure of the molecule in a graph

Employ DGNs to predict the mapping entropy





Errica, Giulini, RM et al., Front. Mol. Biosci. 2021



$$\mathbf{h}_{v}^{\ell+1} = MLP^{\ell} \Big( \left( 1 + \epsilon^{\ell} \right) * \mathbf{h}_{v}^{\ell} + \sum_{u \in \mathcal{N}_{v}} \mathbf{h}_{u}^{\ell} * e^{\ell} \Big) \Big)$$







#### Tamapin















Training the DGN on the Smap of - Random CG representations - Maximally informative CG representations: vanishing statistical weight









Training the DGN on the Smap of - Random CG representations - Maximally informative CG representations: vanishing statistical weight

The DGN succeds in predicting the S<sub>map</sub> of both classes. **Speedup wrt standard algorithm:** 10<sup>3</sup> (CPU) - 10<sup>5</sup> (GPU)















### Speedup generated by DGNs: quasi-exhaustive exploration of the Smap landscape











Combine the trained DGNs with the Wang-Landau sampling algorithm to reconstruct the density of states  $\Omega(S_{map})$ 



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021











Combine the trained DGNs with the Wang-Landau sampling algorithm to reconstruct the density of states  $\Omega(S_{map})$ 



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021











Combine the trained DGNs with the Wang-Landau sampling algorithm to reconstruct the density of states  $\Omega(S_{map})$ 



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021













Combine the trained DGNs with the Wang-Landau sampling algorithm to reconstruct the density of states  $\Omega(S_{map})$ 



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021











Combine the trained DGNs with the Wang-Landau sampling algorithm to reconstruct the density of states  $\Omega(S_{map})$ 



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021



# **Smap via Deep Graph Networks**



DGNs succeed in capturing the correct statistical weight of CG representations: do not overfit the training set









# **Application to neural systems**



#### Being based on information-theoretic quantities, the mapping entropy protocol is extremely general







# **Application to neural systems**



#### Being based on information-theoretic quantities, the mapping entropy protocol is extremely general







19 (TIFPA)





Hopfield, PNAS (1982) Amit, *Modeling Brain Function: The World of Attractor Neural Networks* (1989) Aldrigo, RM, Potestio (in preparation)









Hopfield, PNAS (1982) Amit, *Modeling Brain Function: The World of Attractor Neural Networks* (1989) Aldrigo, RM, Potestio (in preparation)



#### Nonlinear evolution of the configurations Model's dynamics retrieves memory patterns

Overlaps 
$$m^{\mu} = \frac{1}{N} \sum_{i} \xi_{i}^{\mu} \sigma_{i}$$







Hopfield, PNAS (1982) Amit, *Modeling Brain Function: The World of Attractor Neural Networks* (1989) Aldrigo, RM, Potestio (in preparation)











Hopfield, PNAS (1982) Amit, *Modeling Brain Function: The World of Attractor Neural Networks* (1989) Aldrigo, RM, Potestio (in preparation)



Looking for the **most informative** neurons!




# **Detection of maximally-informative neurons**

### CG'ing the Hopfield model

### 1) Simulate the Hopfield model









# **Detection of maximally-informative neurons**

### CG'ing the Hopfield model









# **Detection of maximally-informative neurons**

### CG'ing the Hopfield model



Aldrigo, RM, Potestio (in preparation)



#### 3) Optimize the selection to detect maximally-informative neurons





- Simulate the high-resolution Hopfield model - **Empirical** reference probability  $p(\sigma_1, ..., \sigma_N)$ 









- Simulate the high-resolution Hopfield model - **Empirical** reference probability  $p(\sigma_1, ..., \sigma_N)$ 



- Select n<N retained neurons  $S_i$
- **Empirical** CG probability  $P_M(S_1, ..., S_n)$
- Empirical backmapped probability  $\bar{p}_M(\sigma_1,...,\sigma_N)$







- Simulate the high-resolution Hopfield model - **Empirical** reference probability  $p(\sigma_1, ..., \sigma_N)$ 



- Select n<N retained neurons  $S_i$
- Empirical CG probability  $P_M(S_1,...,S_n)$
- Empirical backmapped probability  $\bar{p}_M(\sigma_1,...,\sigma_N)$

Aldrigo, RM, Potestio (in preparation)



Resolution of the neuron selection

$$\mathcal{H}_M = -\sum_{\{S_i\}} P(S_1, ..., S_n) \ln P(S_1, ..., S_n)$$

- Depends on the specific selection
- **Decreases** by decreasing the number of retained neurons





- Simulate the high-resolution Hopfield model - **Empirical** reference probability  $p(\sigma_1, ..., \sigma_N)$ 



- Select n<N retained neurons  $S_i$
- **Empirical** CG probability  $P_M(S_1, ..., S_n)$
- Empirical backmapped probability  $\bar{p}_M(\sigma_1,...,\sigma_N)$

Aldrigo, RM, Potestio (in preparation)



**Resolution** of the neuron selection

$$\mathcal{H}_M = -\sum_{\{S_i\}} P(S_1, ..., S_n) \ln P(S_1, ..., S_n)$$

- **Depends** on the specific selection

- **Decreases** by decreasing the number of retained neurons

**Information loss** generated by the selection: mapping entropy

$$S_M^{map} = \sum_{\{\sigma_i\}} p(\{\sigma_i\}) \ln\left(\frac{p(\{\sigma_i\})}{\bar{p}_M(\{\sigma_i\})}\right)$$

- **Depends** on the specific selection

- **Increases** by decreasing the number of retained neurons





- Simulate the high-resolution Hopfield model - **Empirical** reference probability  $p(\sigma_1, ..., \sigma_N)$ 



- Select n<N retained neurons  $S_i$
- **Empirical** CG probability  $P_M(S_1, ..., S_n)$
- Empirical backmapped probability  $\bar{p}_M(\sigma_1,...,\sigma_N)$

Aldrigo, RM, Potestio (in preparation)



**Resolution** of the neuron selection

$$\mathcal{H}_M = -\sum_{\{S_i\}} P(S_1, ..., S_n) \ln P(S_1, ..., S_n)$$

- **Depends** on the specific selection

- **Decreases** by decreasing the number of retained neurons

Maximally informative neurons: minimize the mapping entropy in the space of possible selections!

**uss** generated by the selection: mapping entropy

$$S_M^{map} = \sum_{\{\sigma_i\}} p(\{\sigma_i\}) \ln\left(\frac{p(\{\sigma_i\})}{\bar{p}_M(\{\sigma_i\})}\right)$$

- **Depends** on the specific selection

- **Increases** by decreasing the number of retained neurons





Maximally informative selection of neurons that **minimize the mapping entropy** Hopfield model with N=100 neurons and 5 memory patterns









Maximally informative selection of neurons that **minimize the mapping entropy** Hopfield model with N=100 neurons and 5 memory patterns









Maximally informative selection of neurons that **minimize the mapping entropy** Hopfield model with N=100 neurons and 5 memory patterns













#### Aldrigo, RM, Potestio (in preparation)



### Maximally informative selection of neurons that **minimize the mapping entropy**







2.0



Strongly interacting retained neurons, weakly coupled with the integrated ones

2

Aldrigo, RM, Potestio (in preparation)



### Maximally informative selection of neurons that **minimize the mapping entropy**









2.0



Strongly interacting retained neurons, weakly coupled with the integrated ones

2

Aldrigo, RM, Potestio (in preparation)



### Maximally informative selection of neurons that **minimize the mapping entropy**









2.0



Strongly interacting retained neurons, weakly coupled with the integrated ones

2

Aldrigo, RM, Potestio (in preparation)



### Maximally informative selection of neurons that **minimize the mapping entropy**















Aldrigo, RM, Potestio (in preparation)



### Maximally informative selection of neurons that **minimize the mapping entropy**









- ●
- space





• CG filters as a mean to gain insight on the behavior of complex systems

CG'ing procedures result in a loss of statistical information

• Maximally-informative CG representations single out biologically relevant residues of a macromolecular system in an unsupervised manner

• Deep Graph networks can be employed to speed up information loss calculations and to achieve a quasi-exhaustive exploration of the mapping

• Application of the mapping entropy protocol to a neural system: maximally informative neurons in a Hopfield model

 Phase separation of the maximally informative representations of the system depending on the observational level of detail

# Thank you for the attention!



### Acknowledgments

Potestio Lab: (& former members)

Raffaello Potestio Thomas Tarenzi, now at UoB (UK) Giovanni Mattiotti Marco Giulini, now at UU (NL) Roberto Menichetti Elio Fiorentini, now at AM (IT) Marta Rigoli, now at CIBIO (IT) Margherita Mele Lorenzo Petrolli Manuel Micheloni Camilla Spreti













Giulini, RM, Shell, Potestio, JCTC 2020







#### Giulini, RM, Shell, Potestio, JCTC 2020









#### Giulini, RM, Shell, Potestio, JCTC 2020



erc









Giulini, RM, Shell, Potestio, JCTC 2020











Optimized solutions are isolated points in the mapping entropy landscape

Atoms conserved with high probability are likely present in many uncorrelated solutions

Giulini, RM, Shell, Potestio, JCTC 2020















erc







erc

Search for mappings that minimize the information loss within each basin







erc



•











•

















erc



















erc









erc



**Biologically relevant residues slightly** "deactivate" between Apo and Holo











erc





**Biologically relevant residues slightly** "deactivate" between Apo and Holo







# **Application: G6 Antibody and VEGF-A antigen**



Galante, Tarenzi, Potestio, preliminary results





### **Application: G6 Antibody and VEGF-A antigen**



Galante, Tarenzi, Potestio, preliminary results



erc






Galante, Tarenzi, Potestio, preliminary results



erc





erc

optimizations in each state





erc







erc



#### The curse of the space dimensionality

An exhaustive sampling of all mappings is unfeasible











erc



#### The curse of the space dimensionality

An exhaustive sampling of all mappings is unfeasible



Explore according to the distance  $\mathcal{D}(M, M')$ 













#### Speedup generated by DGNs: quasi-exhaustive exploration of the Smap landscape

RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021



### Smap via deep graph networks



## S<sub>map</sub> via deep graph networks









#### erc

#### Speedup generated by DGNs: <u>quasi-exhaustive</u> exploration of the S<sub>map</sub> landscape

Combine the trained DGN with the Wang-Landau sampling algorithm to reconstruct the density of states of S<sub>map</sub>



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021



### S<sub>map</sub> via deep graph networks









Combine the trained DGN with the Wang-Landau sampling algorithm to reconstruct the density of states of Smap



RM, Giulini, Potestio, EPJB 2021 Errica, Giulini, RM et al., Front. Mol. Biosci. 2021



erc

### S<sub>map</sub> via deep graph networks









Combine the trained DGN with the Wang-Landau sampling algorithm to reconstruct the density of states of S<sub>map</sub>



Errica, Giulini, RM et al., Front. Mol. Biosci. 2021



erc

## **S**<sub>map</sub> via deep graph networks









Combine the trained DGN with the Wang-Landau sampling algorithm to reconstruct the density of states of S<sub>map</sub>





erc

## Smap via deep graph networks



## Moving towards the SG phase

N=100 neurons, p= 4 patterns







erc

N=100 neurons, p= 5 patterns











erc

#### Moving towards the SG phase





#### $S = \left( \begin{array}{c|c} s^* s^* \\ \hline s \ s^* \end{array} \right)$ $\frac{s^*s}{s\,s}$



$$\begin{split} & \text{Semidispersion} \\ & \rho = \frac{\langle |J| \rangle_{s^*s^*} - \langle |J| \rangle_{ss^*}}{\langle |J| \rangle_{s^*s^*} + \langle |J| \rangle_{ss^*}} \end{split}$$











#### DATA != INFORMATION

Simplify "just right" to understand without losing information





erc

#### Simplify to understand





Dimensionality reduction: in terms of what variables?



## Simplify to understand

ATOM	1153	CD	GLN	Х	408
ATOM	1154	OE1	GLN	Х	408
ATOM	1155	NE2	GLN	Х	408
ATOM	1156	HE21	GLN	Х	408
ATOM	1157	HE22	GLN	Х	408
ATOM	1158	С	GLN	Х	408
ATOM	1159	0	GLN	Х	408
ATOM	1160	Ν	GLN	Х	409
ATOM	1161	Η	GLN	Х	409
ATOM	1162	CA	GLN	Х	409
ATOM	1163	HA	GLN	Х	409
ATOM	1164	CB	GLN	Х	409
ATOM	1165	HB1	GLN	Х	409
ATOM	1166	HB2	GLN	Х	409



erc

0
0
0
0
0
0
0
0
0
0
0
0
0
0



# Simplify to understand

<b>Collective variables</b>
-----------------------------

ATOM	1153	CD	GLN	Х	408
ATOM	1154	OE1	GLN	Х	408
ATOM	1155	NE2	GLN	Х	408
ATOM	1156	HE21	GLN	Х	408
ATOM	1157	HE22	GLN	Х	408
ATOM	1158	С	GLN	Х	408
ATOM	1159	0	GLN	Х	408
ATOM	1160	Ν	GLN	Х	409
ATOM	1161	Η	GLN	Х	409
ATOM	1162	CA	GLN	Х	409
ATOM	1163	HA	GLN	Х	409
ATOM	1164	CB	GLN	Х	409
ATOM	1165	HB1	GLN	Х	409
ATOM	1166	HB2	GLN	Х	409



110.020	92.423	42.297	0.00	0.00	
110.283	93.074	43.331	0.00	0.00	
108.983	91.618	42.373	0.00	0.00	
108.656	91.204	41.511	0.00	0.00	
108.471	91.768	43.225	0.00	0.00	
113.167	95.332	40.798	0.00	0.00	
113.764	95.563	39.715	0.00	0.00	
113.578	95.654	42.042	0.00	0.00	
113.018	95.252	42.780	0.00	0.00	
114.881	96.418	42.260	0.00	0.00	
115.146	96.965	41.356	0.00	0.00	
114.444	97.423	43.339	0.00	0.00	
115.314	97.840	43.844	0.00	0.00	
113.900	96.872	44.111	0.00	0.00	



# Simplify to understand

<b>Collective variables</b>
-----------------------------

ATOM	1153	CD	GLN	Х	408
ATOM	1154	OE1	GLN	Х	408
ATOM	1155	NE2	GLN	Х	408
ATOM	1156	HE21	GLN	Х	408
ATOM	1157	HE22	GLN	Х	408
ATOM	1158	С	GLN	Х	408
ATOM	1159	0	GLN	Х	408
ATOM	1160	Ν	GLN	Х	409
ATOM	1161	Η	GLN	Х	409
ATOM	1162	CA	GLN	Х	409
ATOM	1163	HA	GLN	Х	409
ATOM	1164	CB	GLN	Х	409
ATOM	1165	HB1	GLN	Х	409
ATOM	1166	HB2	GLN	Х	409



110.020	92.423	42.297	0.00	0.00	
110.283	93.074	43.331	0.00	0.00	
108.983	91.618	42.373	0.00	0.00	
108.656	91.204	41.511	0.00	0.00	
108.471	91.768	43.225	0.00	0.00	
113.167	95.332	40.798	0.00	0.00	
113.764	95.563	39.715	0.00	0.00	
113.578	95.654	42.042	0.00	0.00	
113.018	95.252	42.780	0.00	0.00	
114.881	96.418	42.260	0.00	0.00	
115.146	96.965	41.356	0.00	0.00	
114.444	97.423	43.339	0.00	0.00	
115.314	97.840	43.844	0.00	0.00	
113.900	96.872	44.111	0.00	0.00	

#### Low-D embedding



Tribello, Gasparotto, Front. Mol. Biosci. 2019





