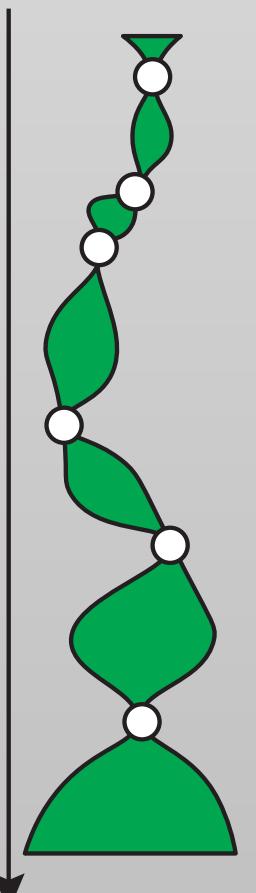
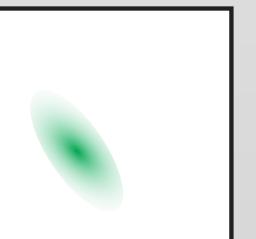
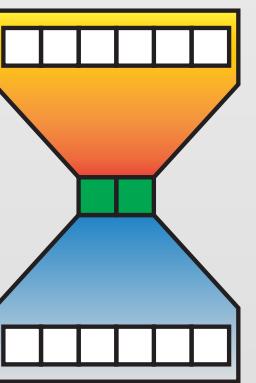
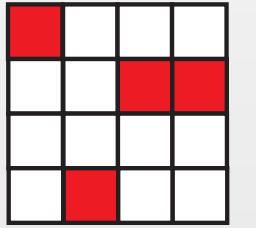
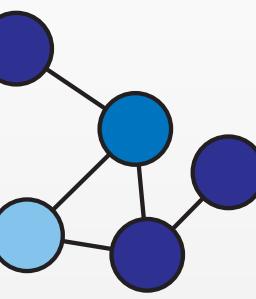
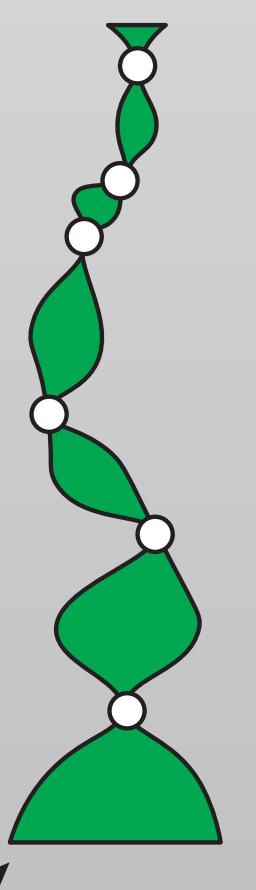
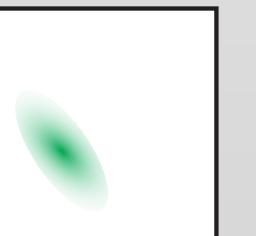
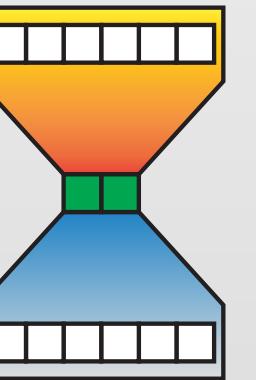
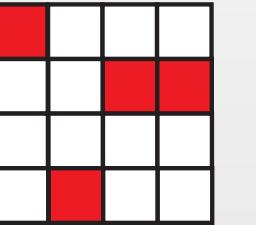
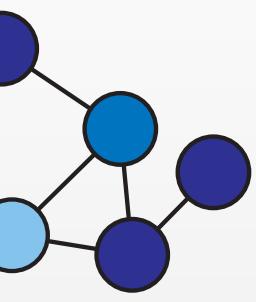


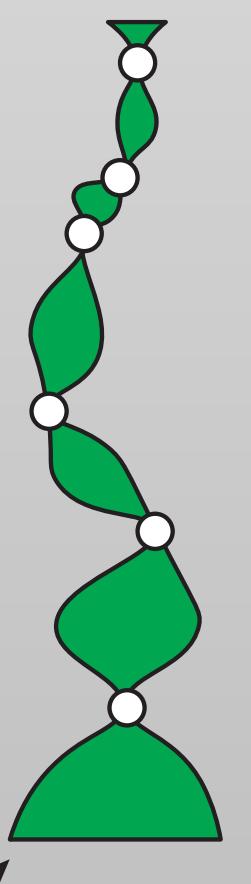
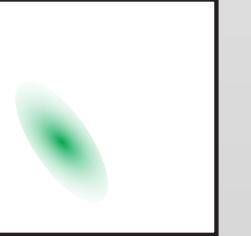
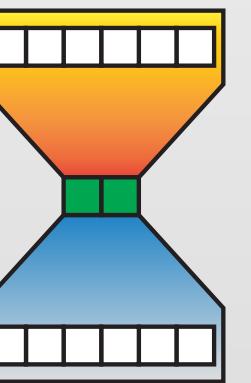
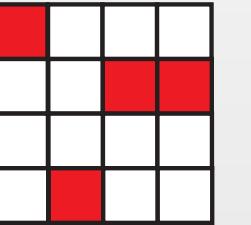
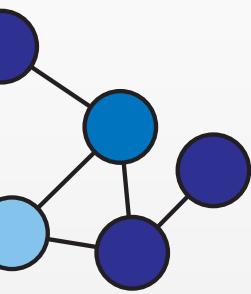
Deep Learning Methods for Learning Phenotype from Genotype

Eric J. Ma
Dept. Biological Engineering, MIT



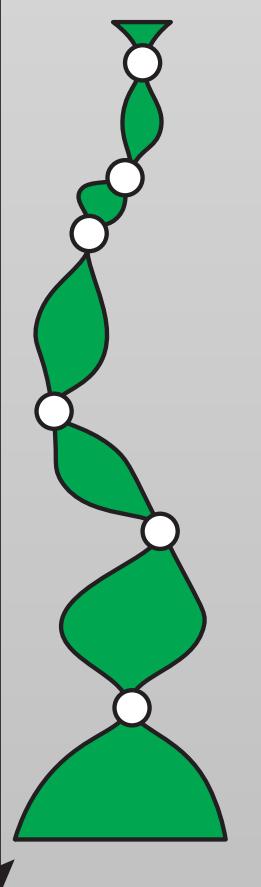
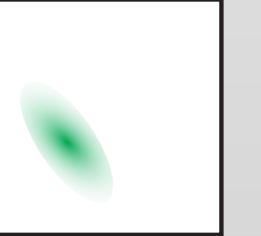
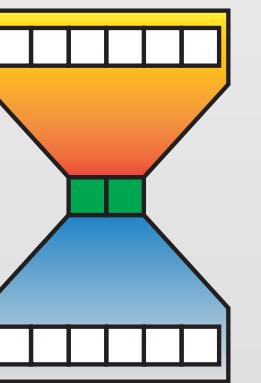
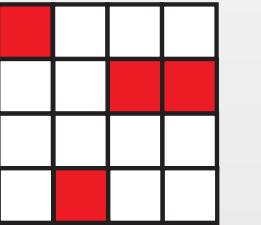
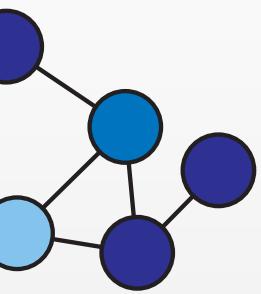
Automatic discovery of structural features predictive of phenotype





Sequence regression ignores positional interactions

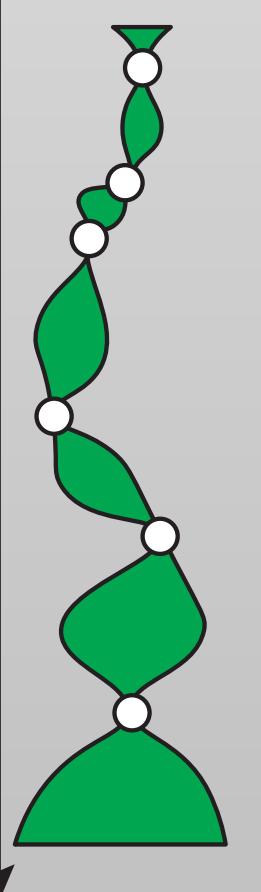
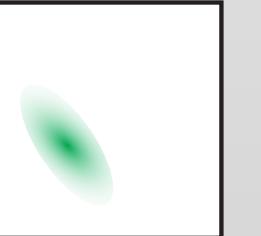
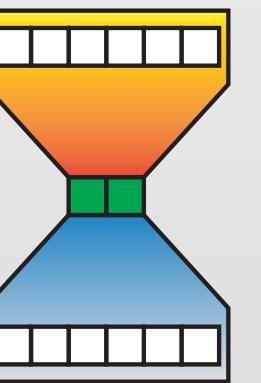
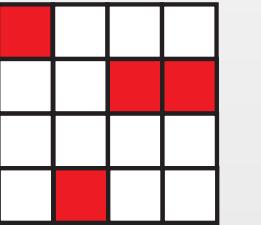
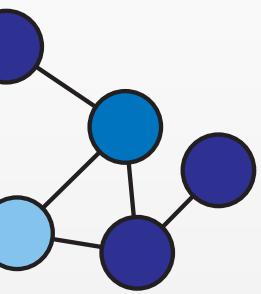
Sequence	DR
PQVTLW <u>Q</u> K PIVTI I KIGG	2 . 4
PQVTLW <u>Q</u> RPIVTI I KIGG	3 . 8
PQVTLW <u>Q</u> RPI <u>L</u> VTI I KIGG	9 . 4
PQVTLW <u>Q</u> RPIVTI I KIGG	3 . 5



Sequence regression ignores positional interactions

Sequence	DR
PQVTLW <u>Q</u> K PIVTI I KIGG	2 . 4
PQVTLW <u>Q</u> RPIVTI I KIGG	3 . 8
PQVTLW <u>Q</u> R <u>P</u> I LVTI I KIGG	9 . 4
PQVTLW <u>Q</u> R <u>P</u> IVTI I KIGG	3 . 5

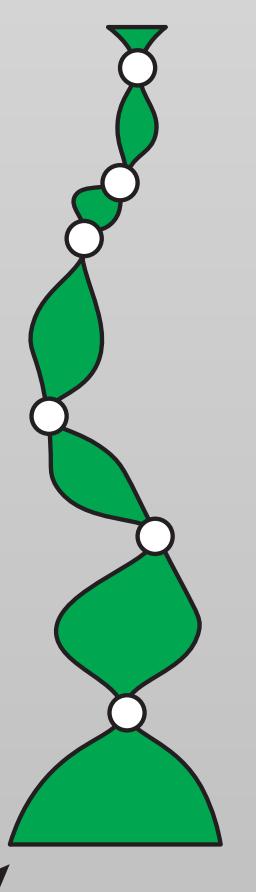
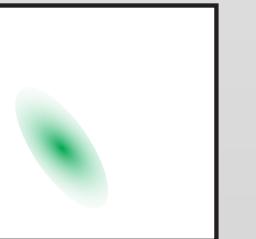
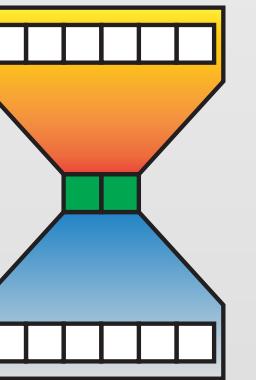
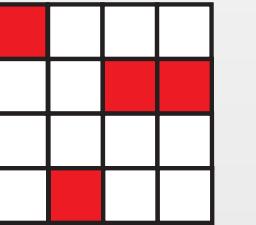
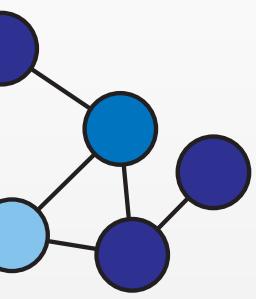




Sequence regression ignores positional interactions

Sequence	DR
PQVTLWQ K P I VTIKIGG	2 . 4
PQVTLWQ R P I VTIKIGG	3 . 8
PQVTLWQ R P L VTIKIGG	9 . 4
PQVTLWQ R P I VTIKIGG	3 . 5

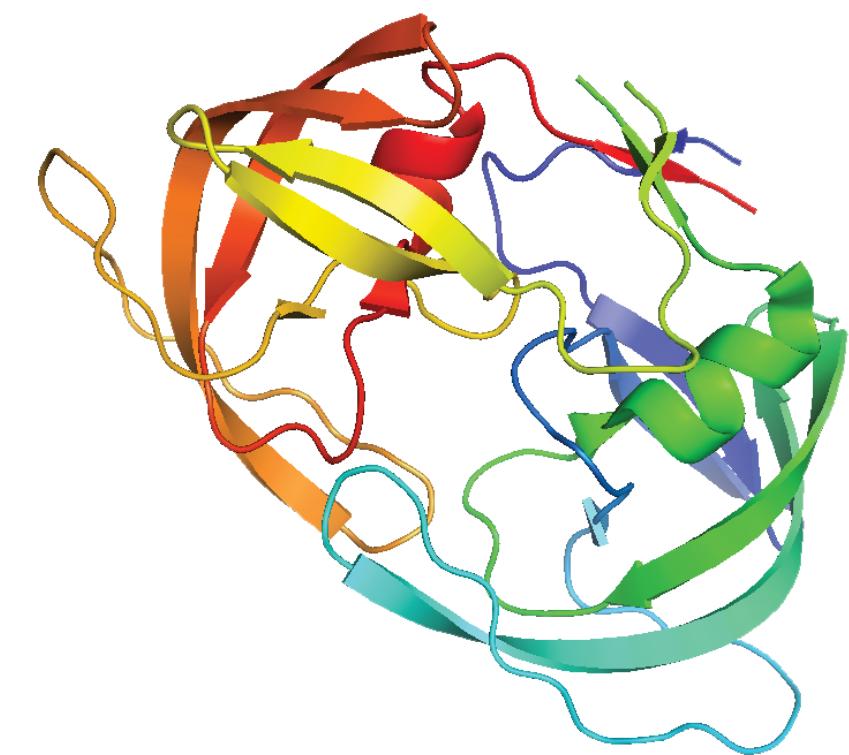
Proteins have a natural graph (network) representation



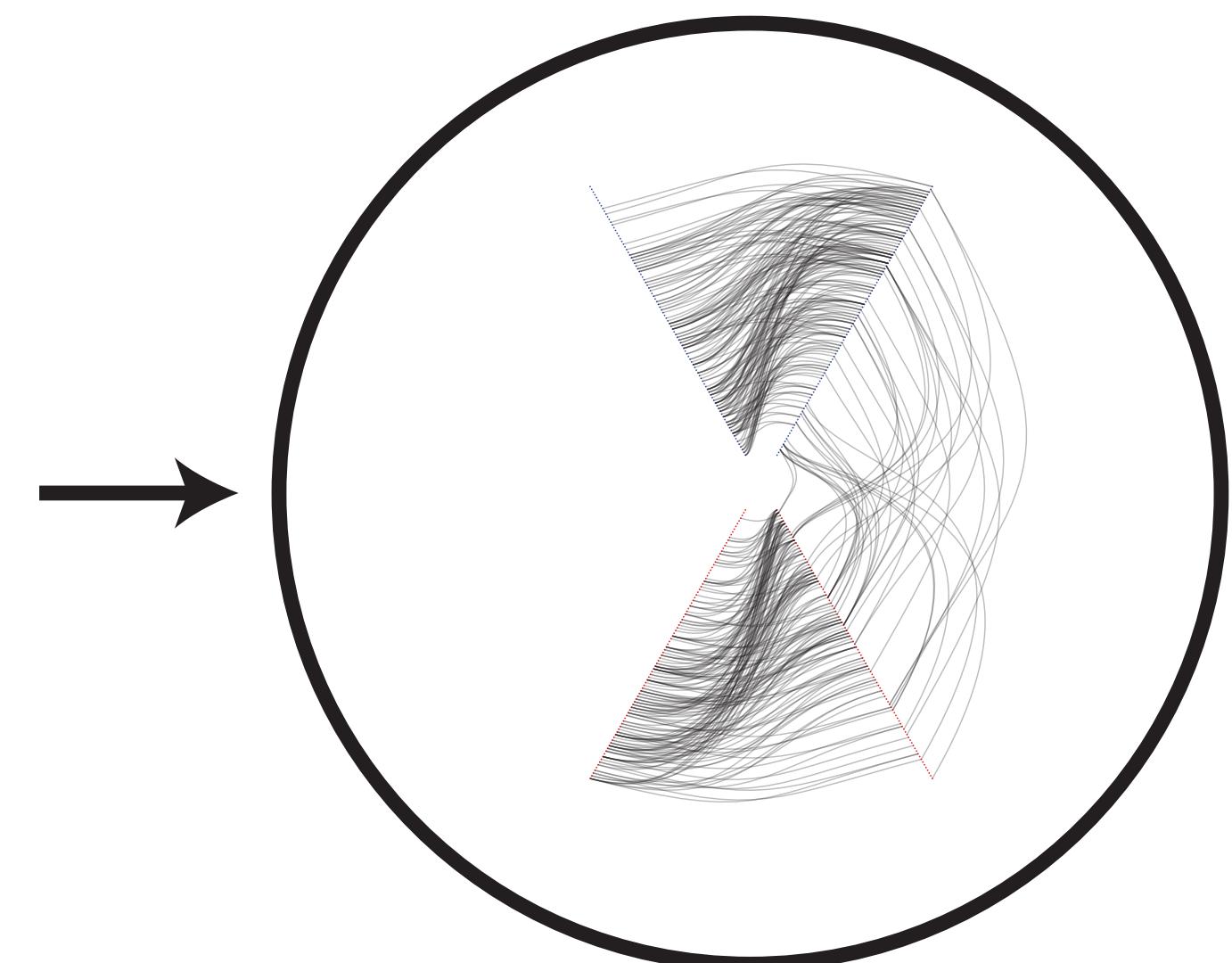
Sequence

PQVTLWQKPIVTIKIGG
PQVTLWQRPIVTIKIGG
PQVTLWQRPLVTIKIGG →
PQVTLWQRPIVTIKIGG

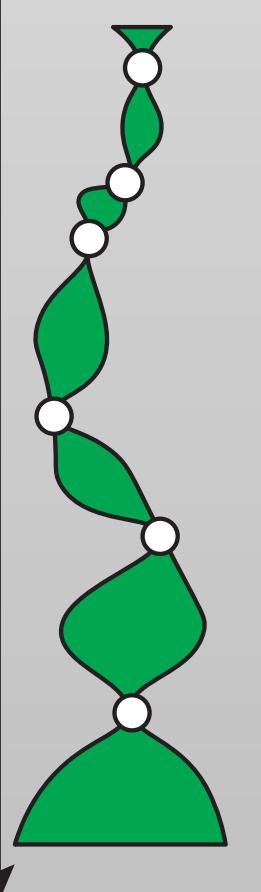
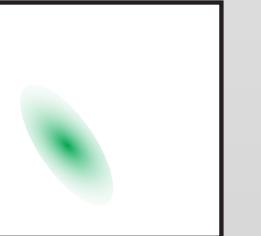
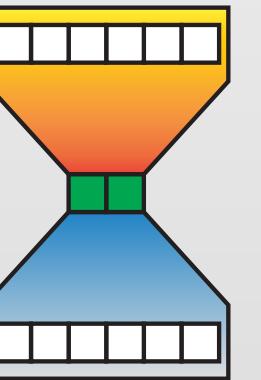
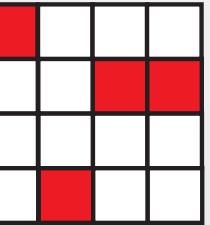
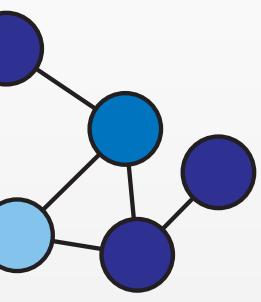
Homology
Model



Network Representation



Proteins have a natural graph (network) representation

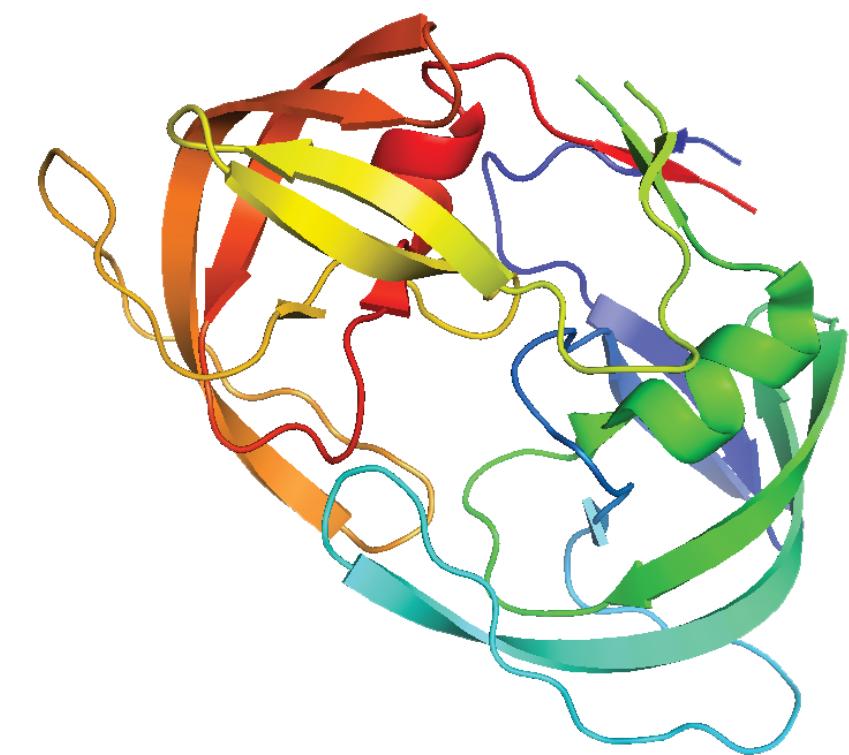


Sequence

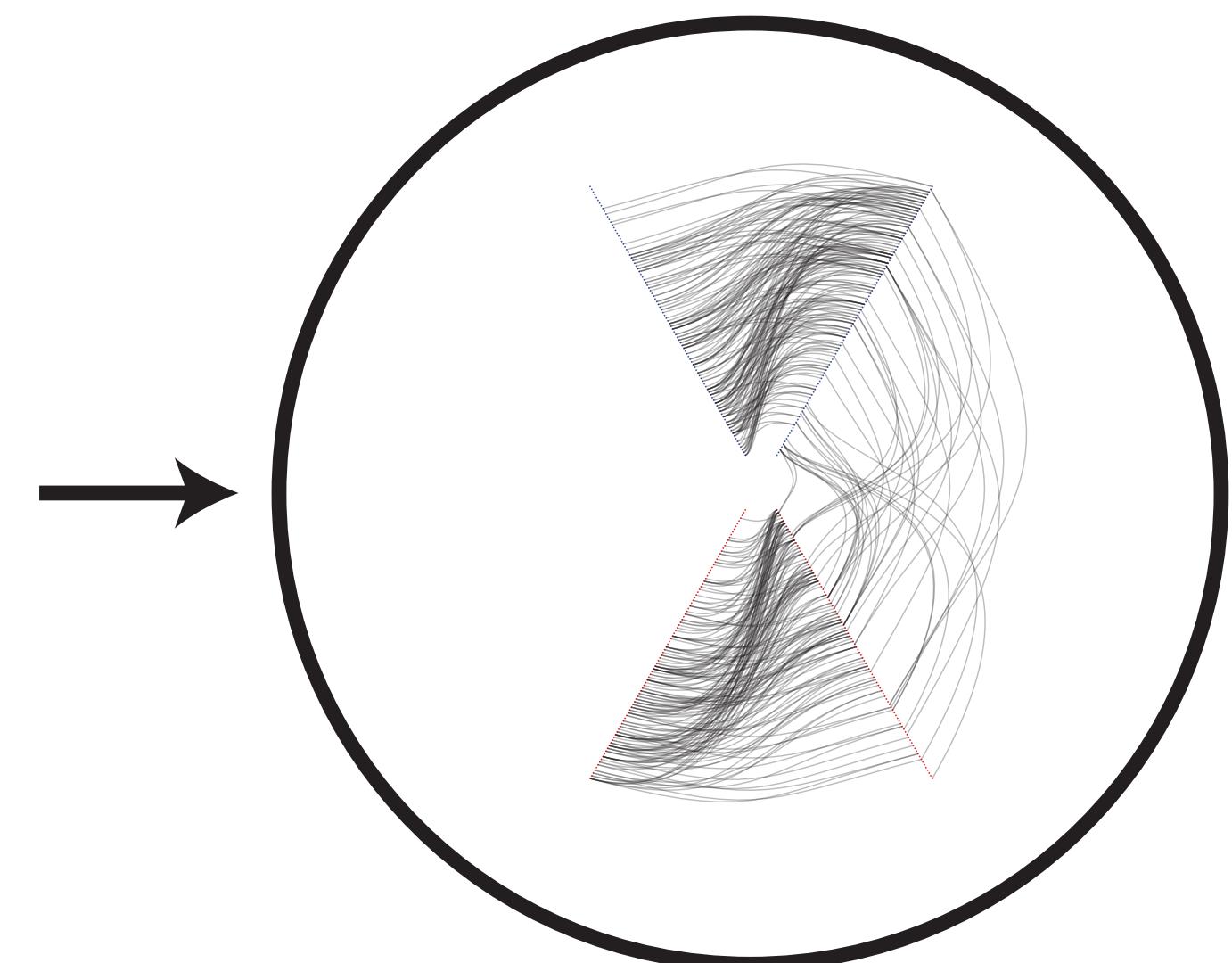
PQVTLWQ**K**PIVTI**K**I**G**
PQVTLWQRPIVTI**K**I**G**
PQVTLWQRPL**V**TI**K**I**G**
PQVTLWQRPIVTI**K**I**G**



Homology
Model



Network Representation

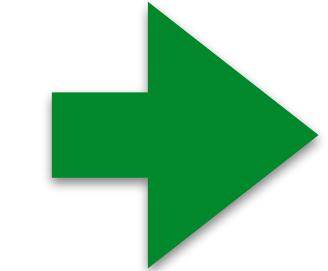


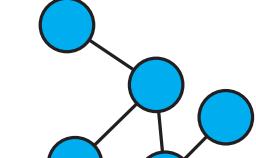
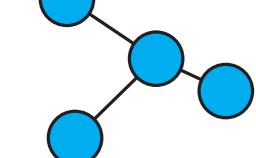
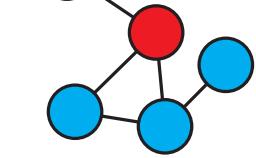
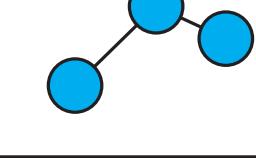
Nodes: Amino acids

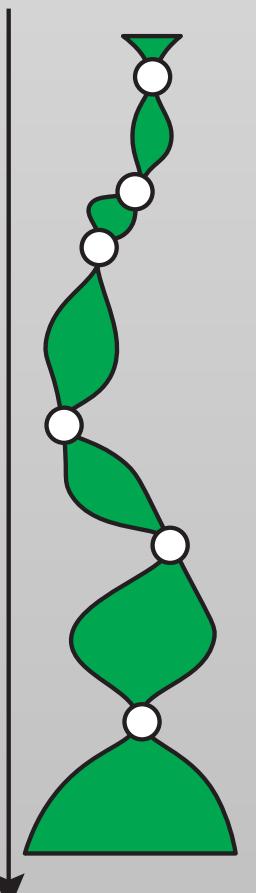
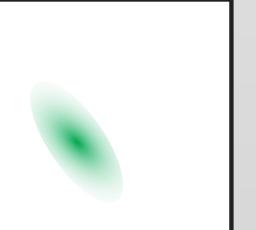
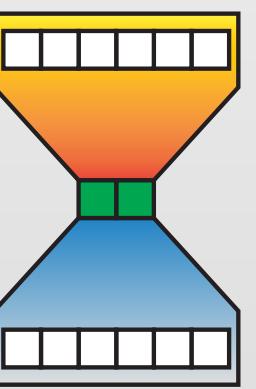
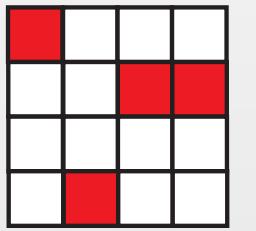
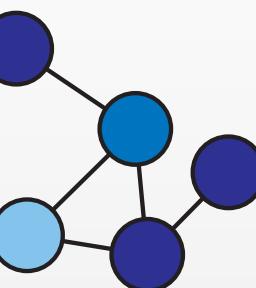
Edges: Biochemical interactions

Can graph regression on proteins help us interpret structure better?

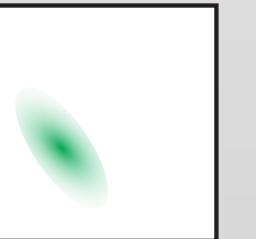
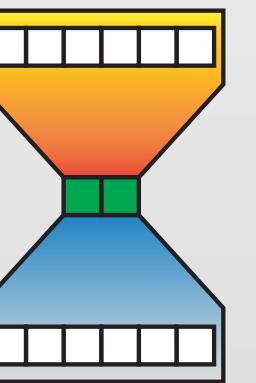
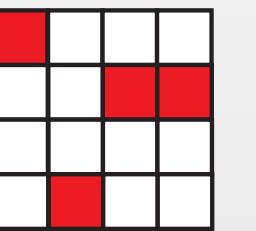
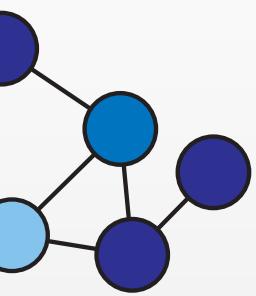
Sequence	DR
PQVTLW <u>Q</u> K PIVTIKIGG	2.4
PQVTLW <u>Q</u> RPIVTIKIGG	3.8
PQVTLW <u>Q</u> R <u>P</u> L VTIKIGG	9.4
PQVTLW <u>Q</u> R <u>P</u> IVTIKIGG	3.5



Structure Graph	DR
	2.4
	3.8
	9.4
	3.5

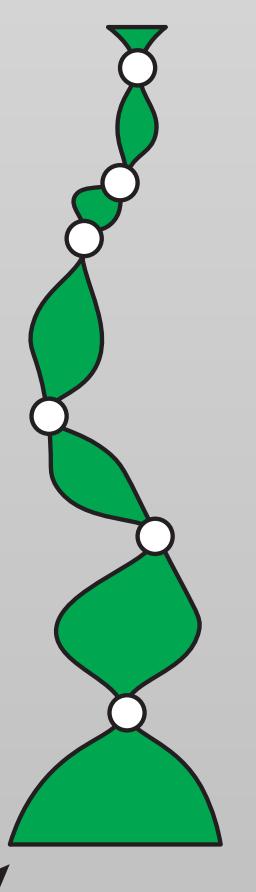
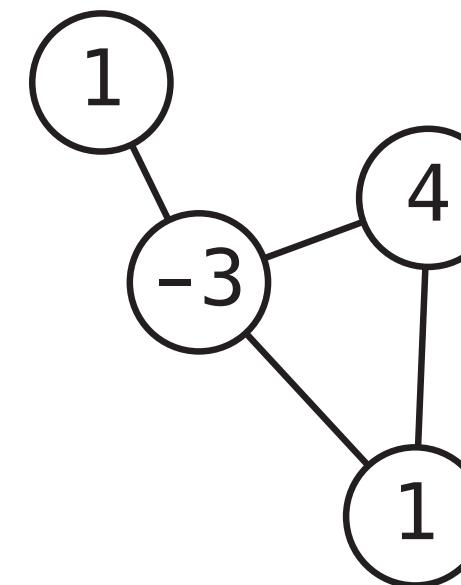


Graph convolutions let us compute unique fingerprints for distance graphs

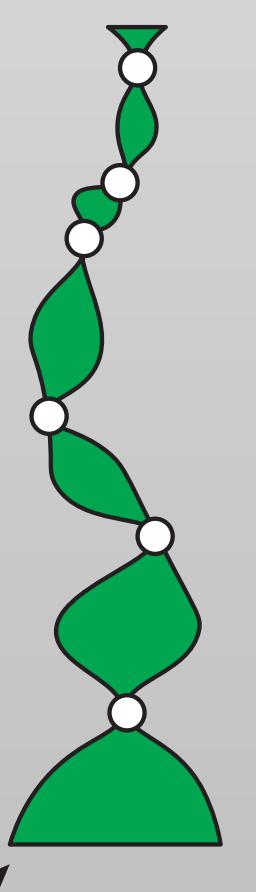
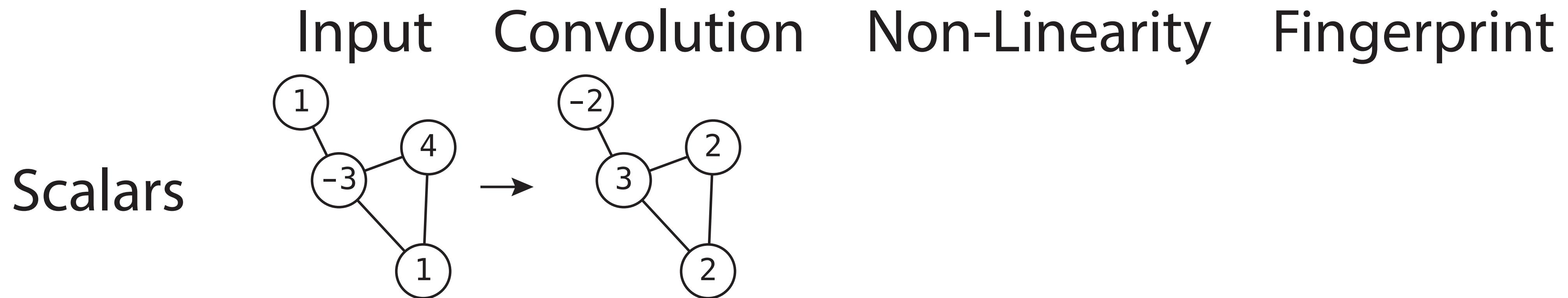
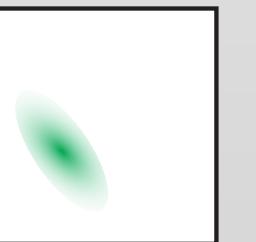
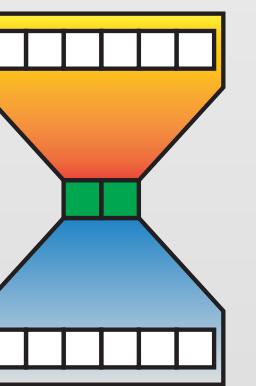
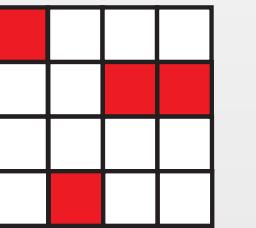
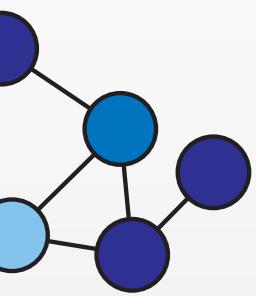


Scalars

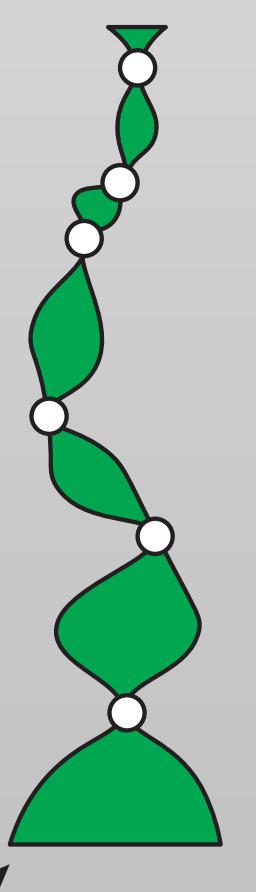
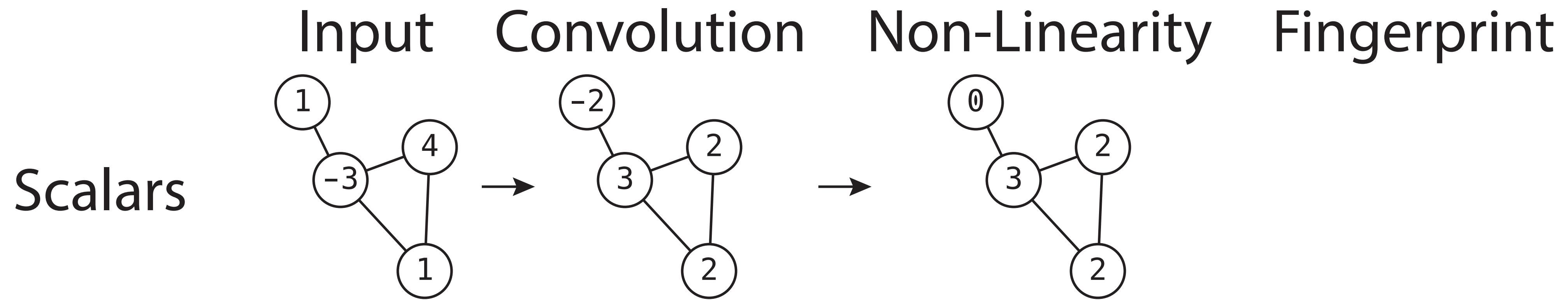
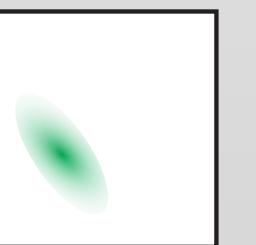
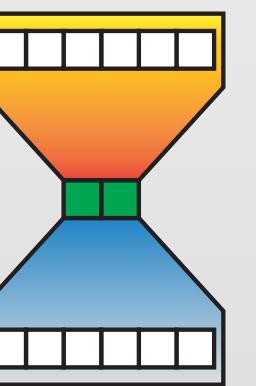
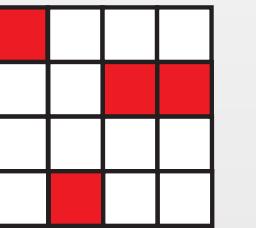
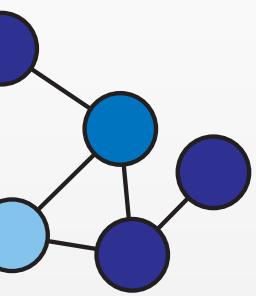
Input Convolution Non-Linearity Fingerprint



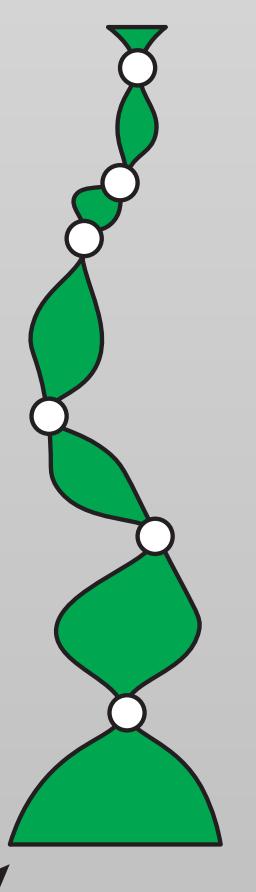
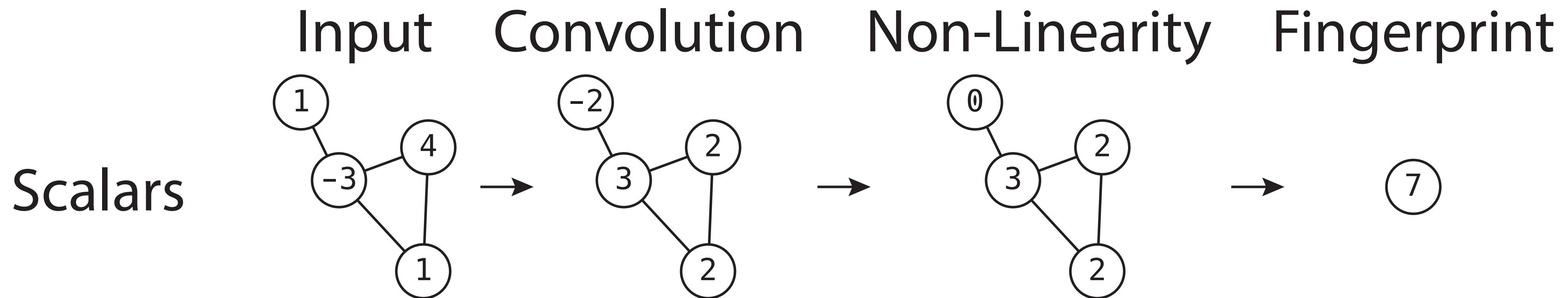
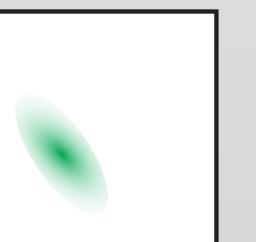
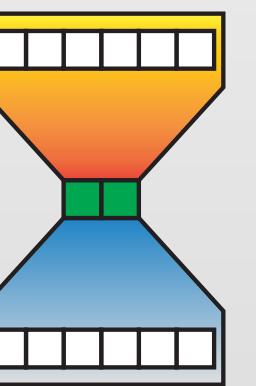
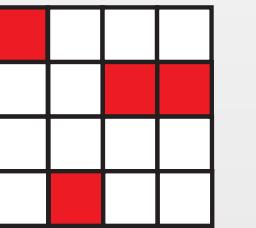
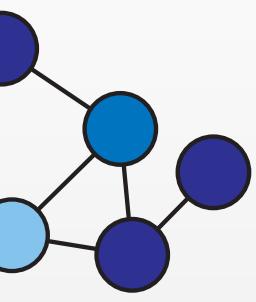
Graph convolutions let us compute unique fingerprints for distance graphs



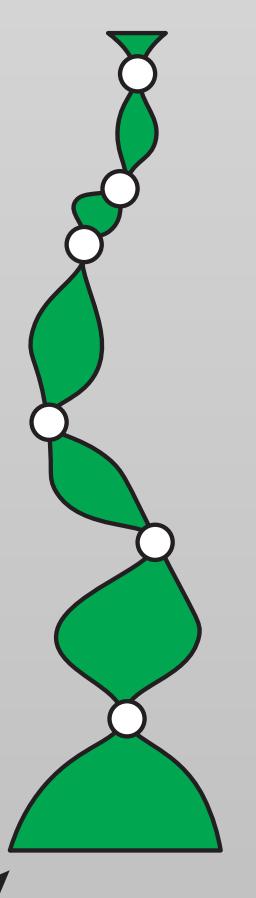
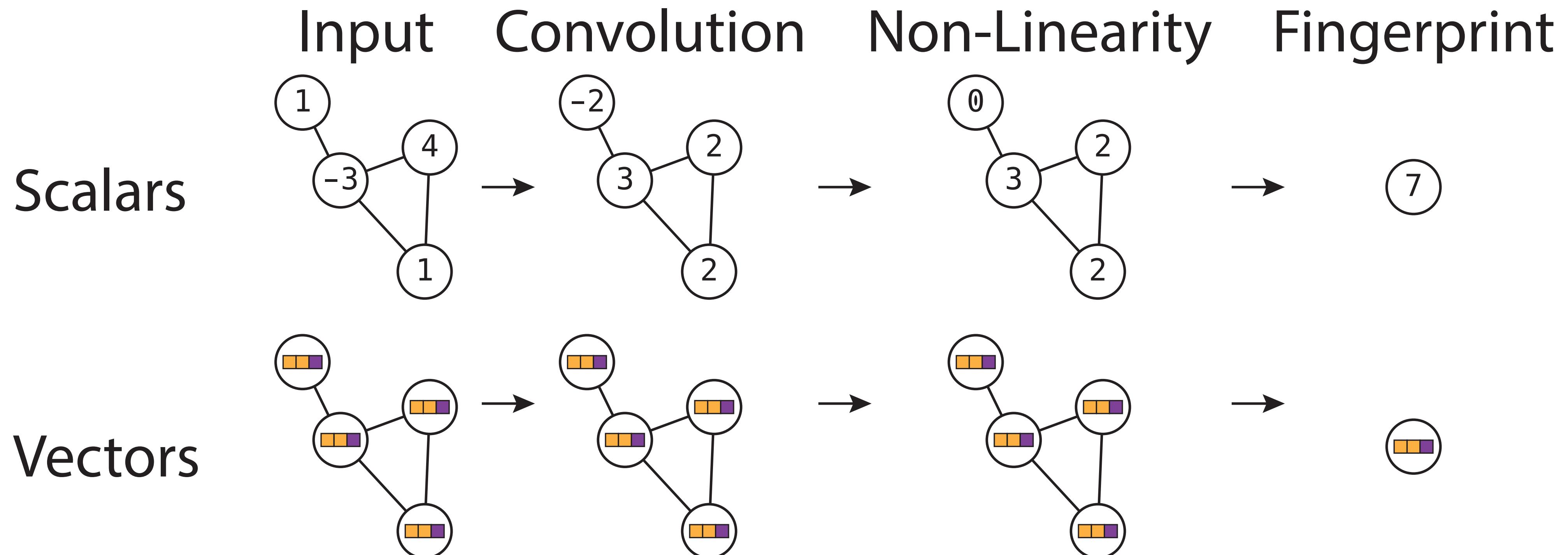
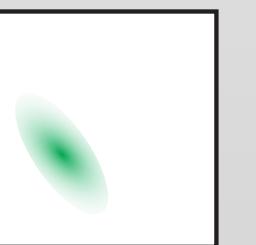
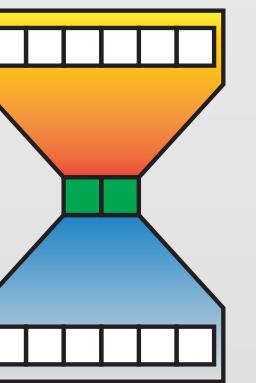
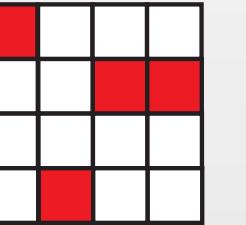
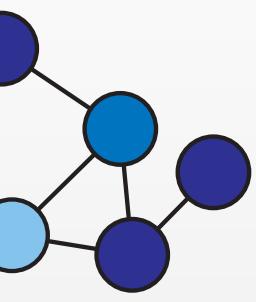
Graph convolutions let us compute unique fingerprints for distance graphs



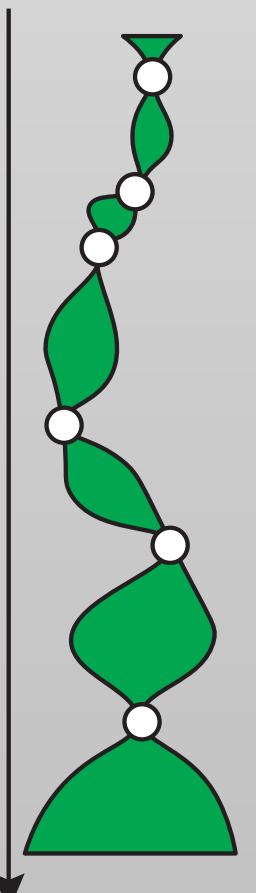
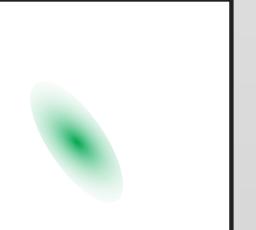
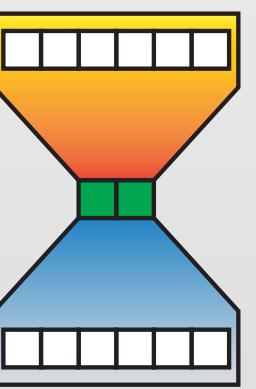
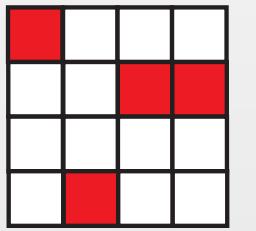
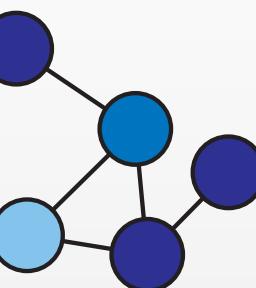
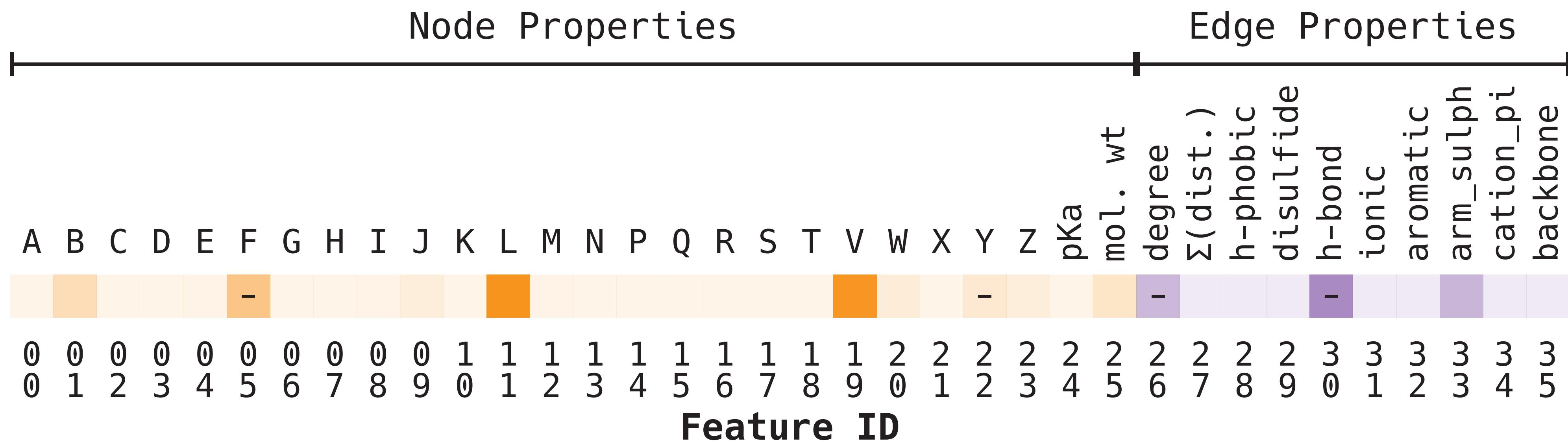
Graph convolutions let us compute unique fingerprints for distance graphs



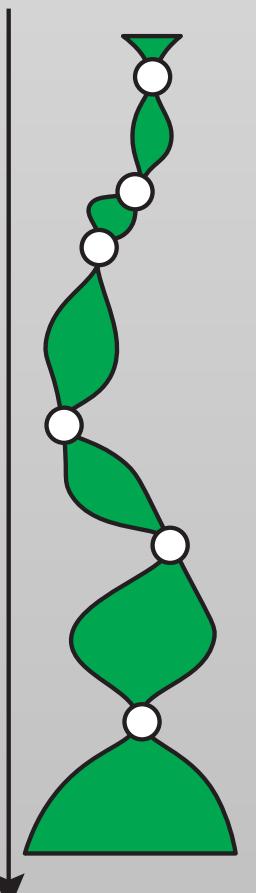
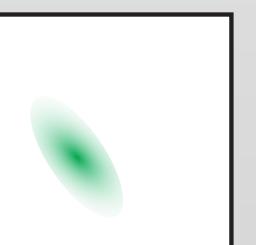
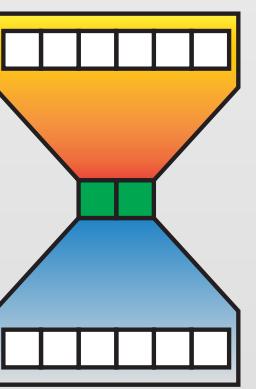
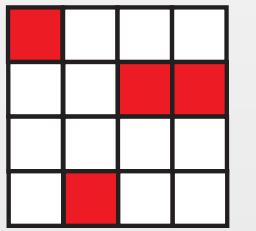
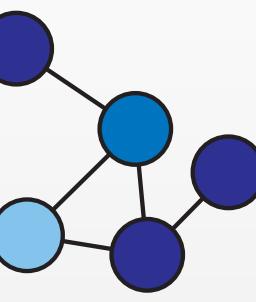
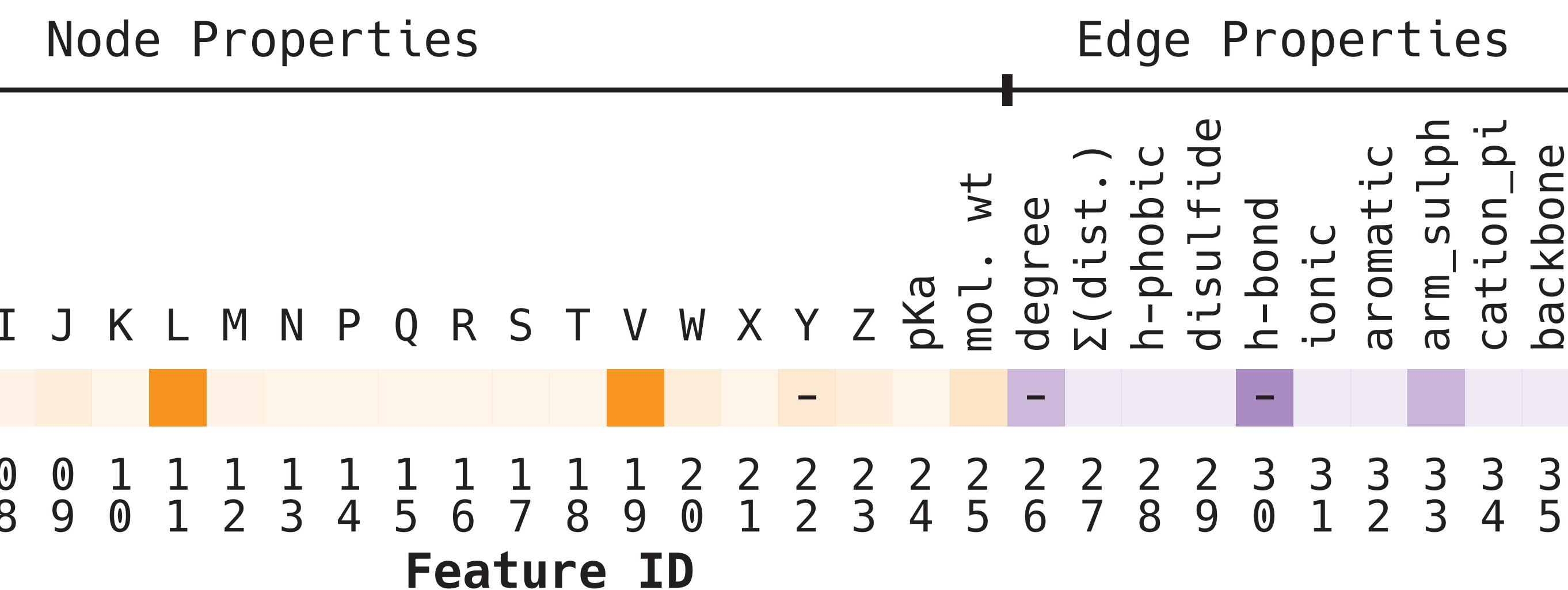
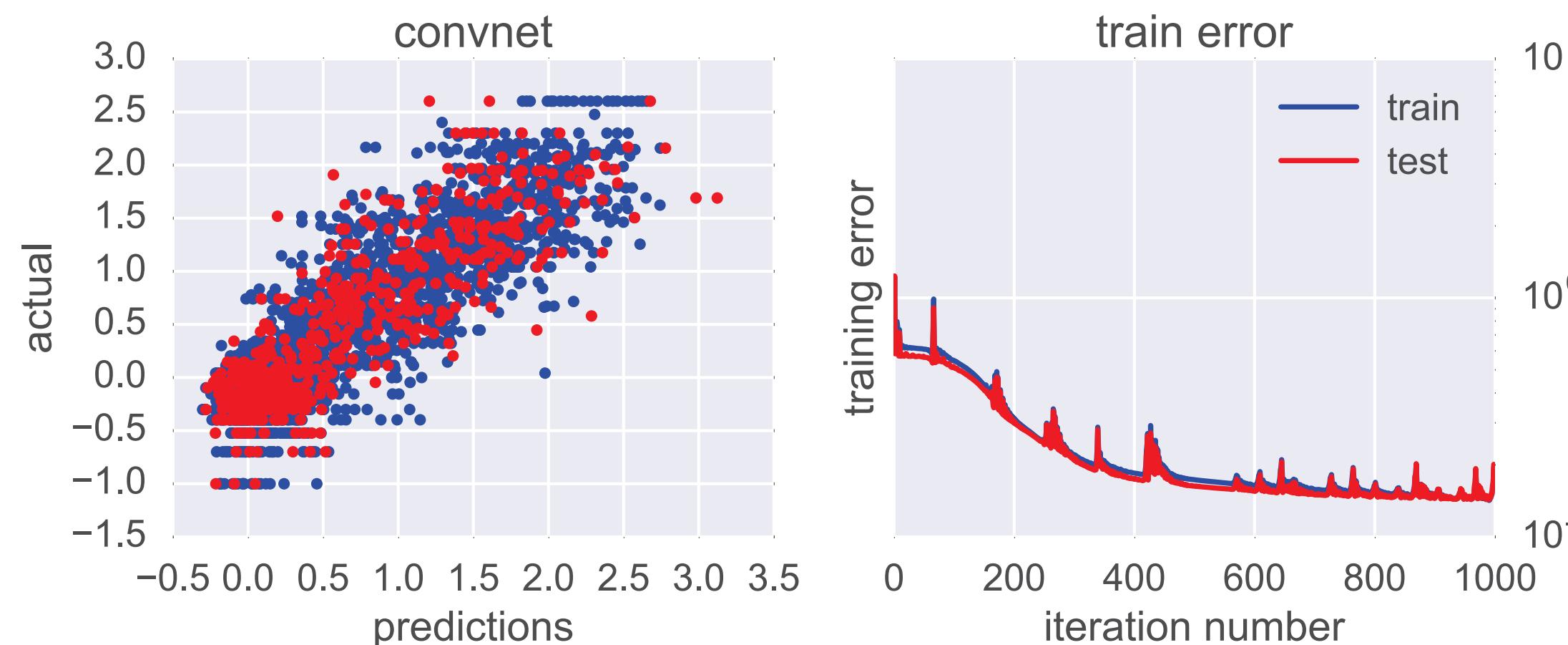
Graph convolutions let us compute unique fingerprints for distance graphs



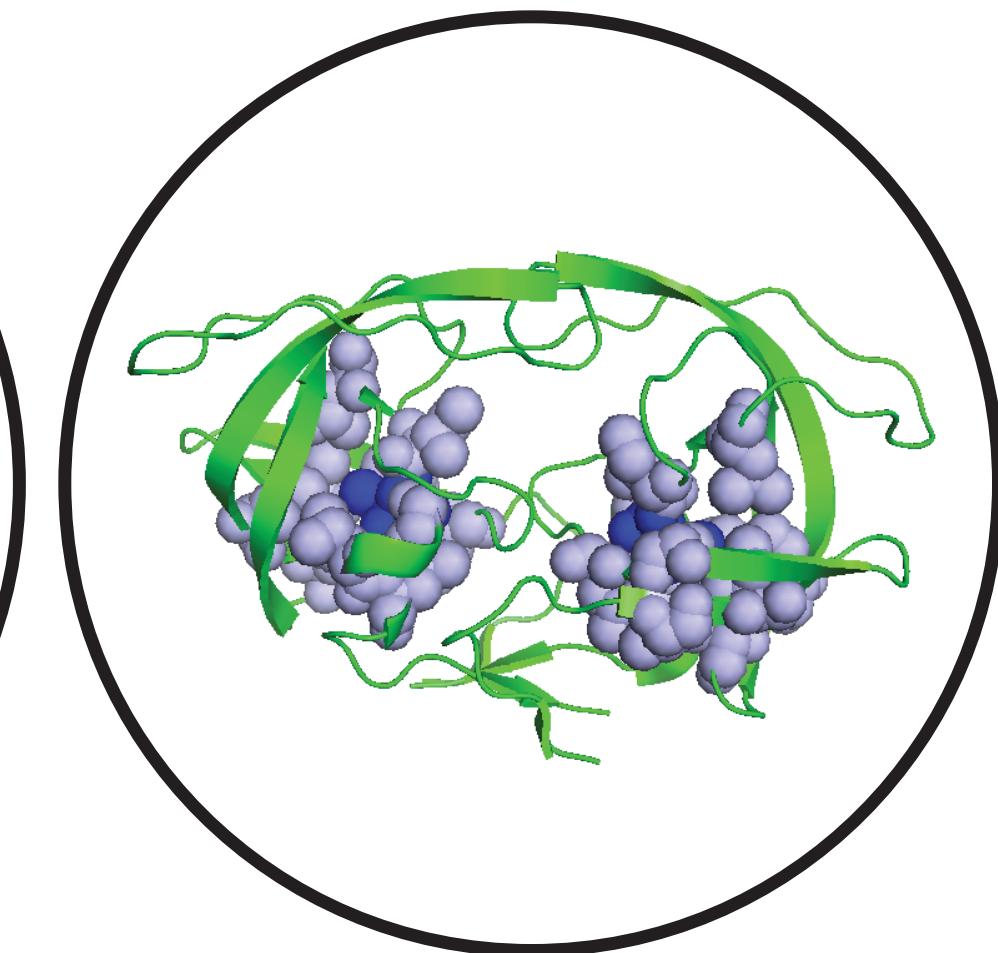
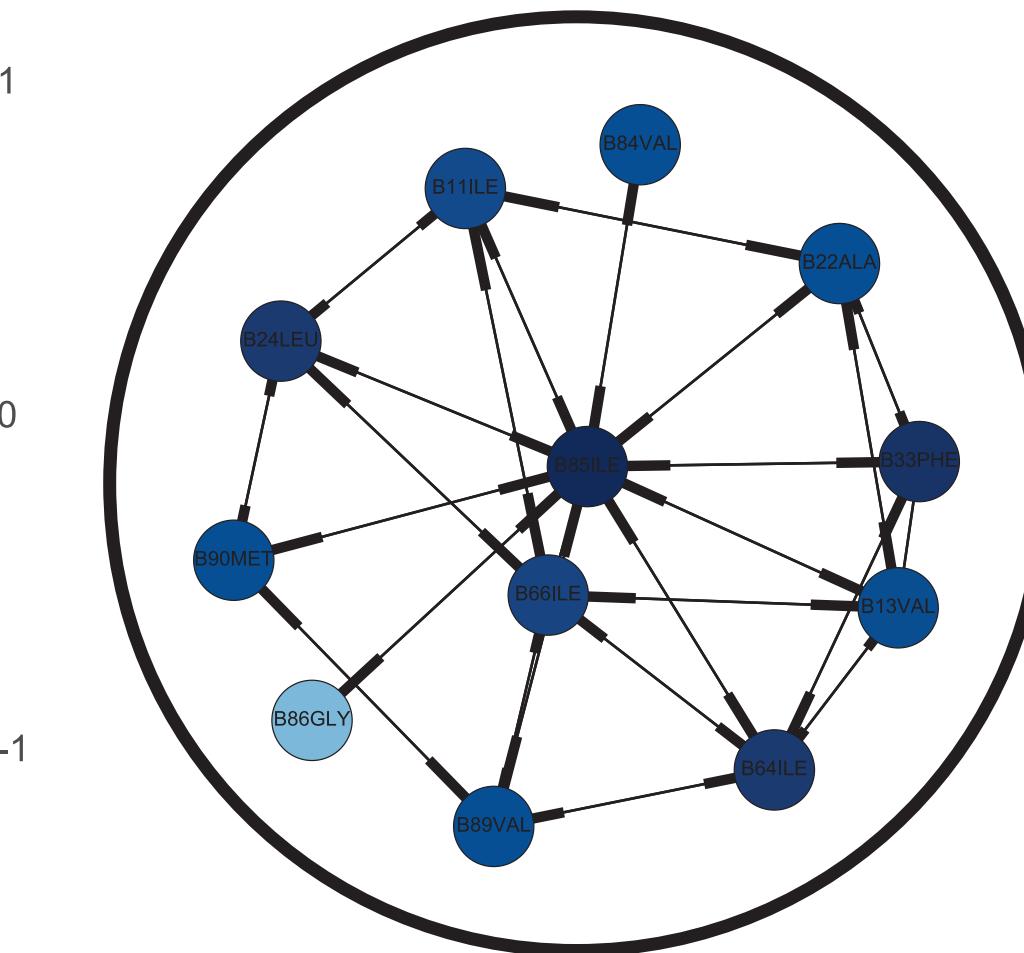
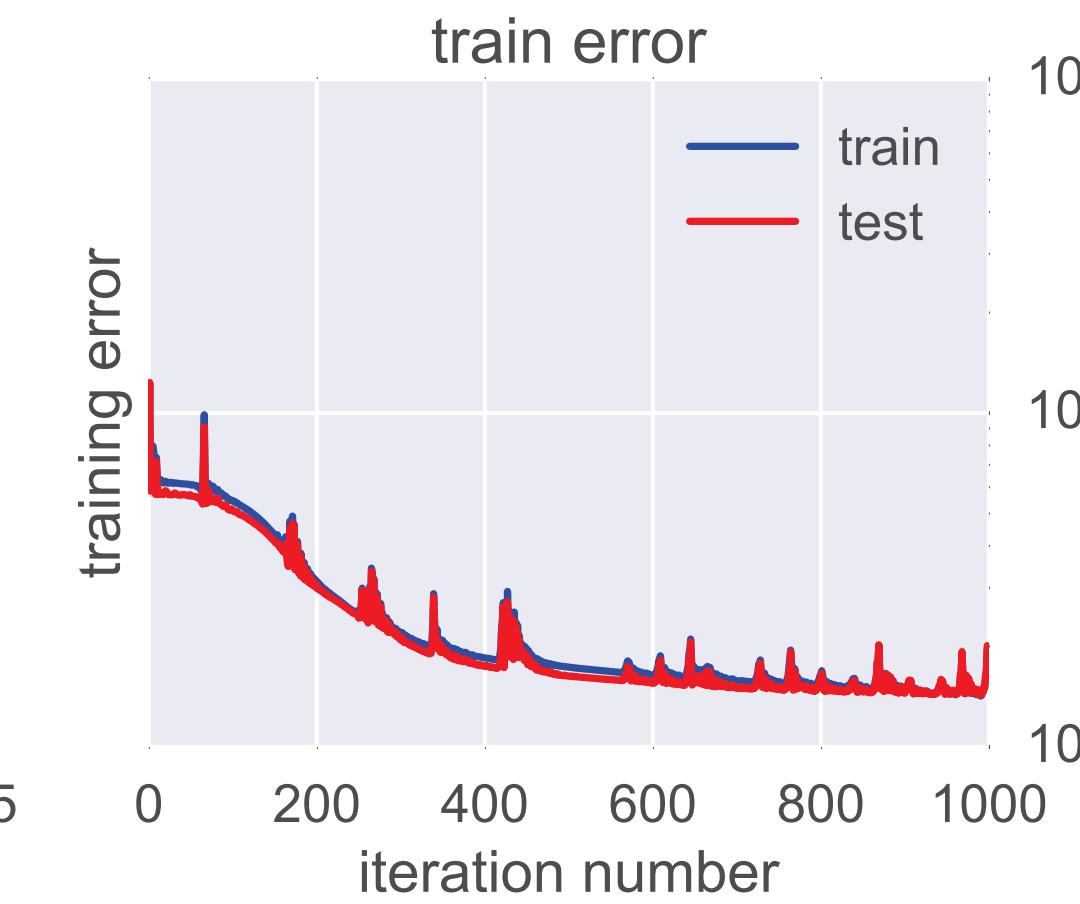
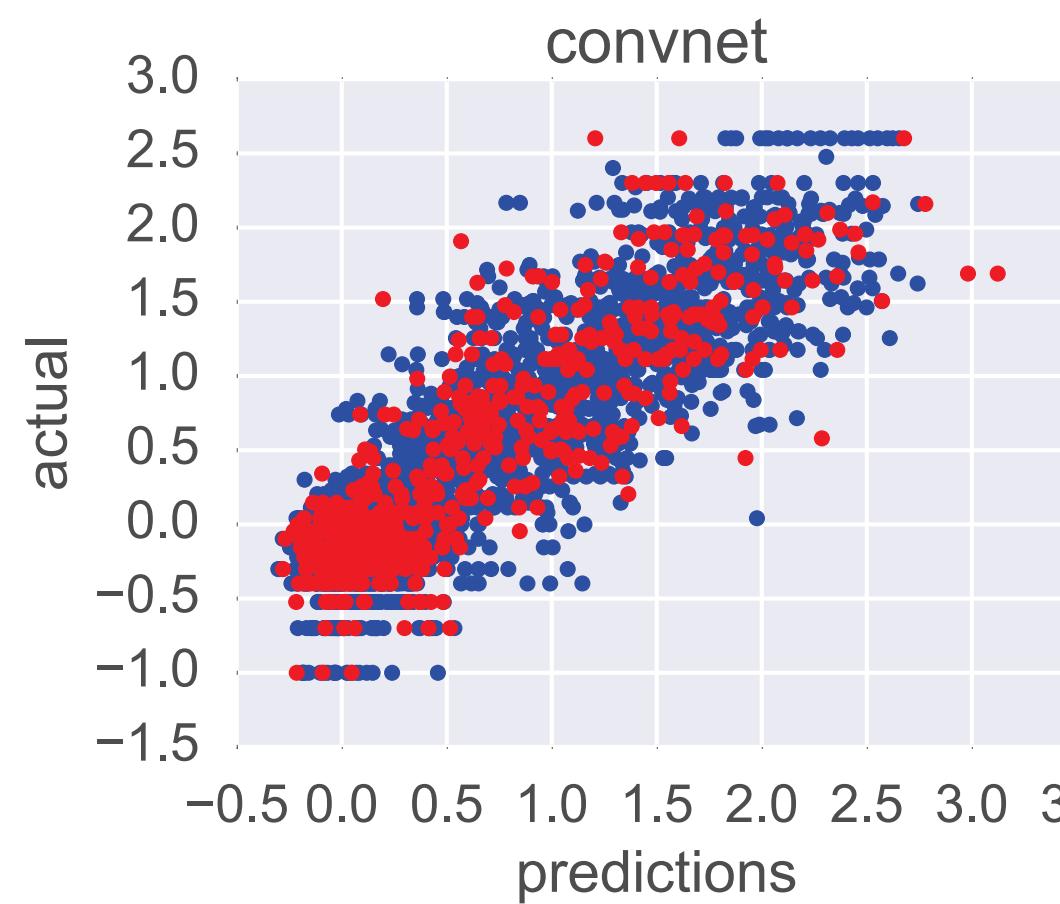
Graph convolutions let us compute unique fingerprints for distance graphs



Graph convolutions automatically learn HIV protease
a.a. neighbourhoods responsible for drug resistance

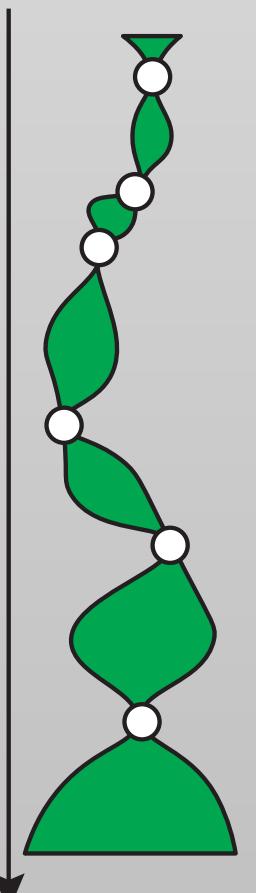
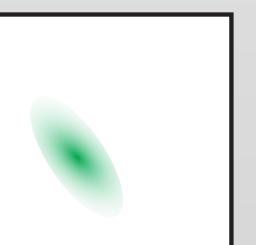
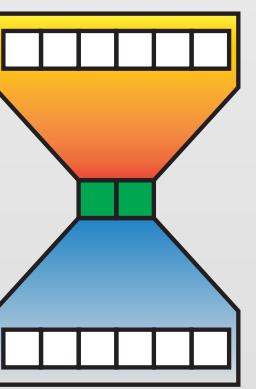
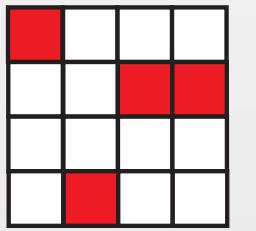
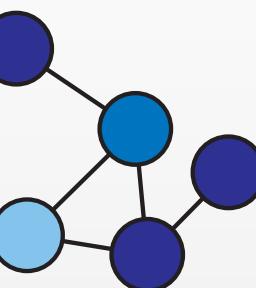
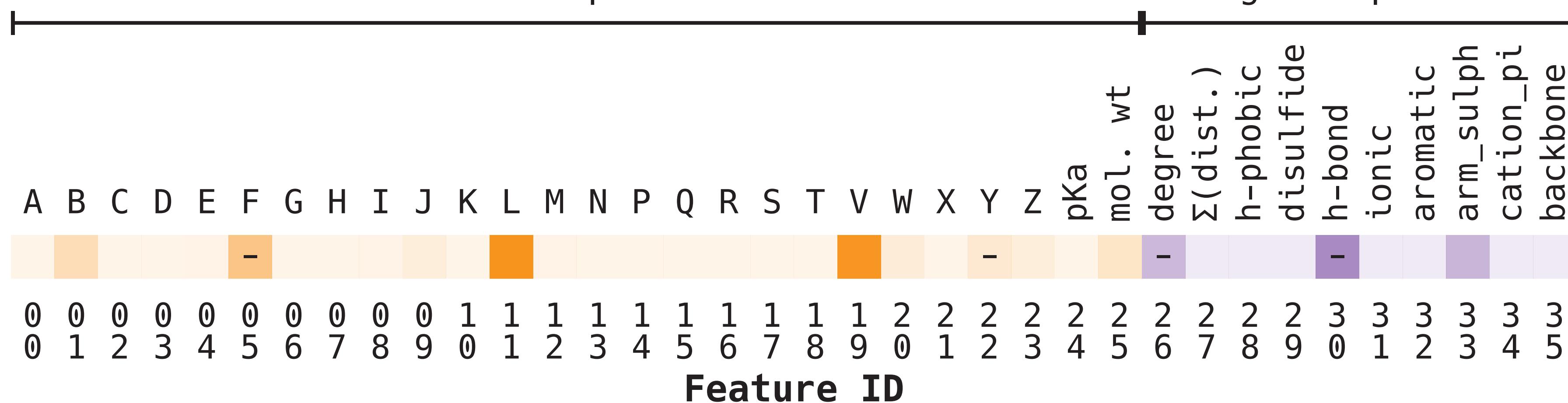


Graph convolutions automatically learn HIV protease a.a. neighbourhoods responsible for drug resistance

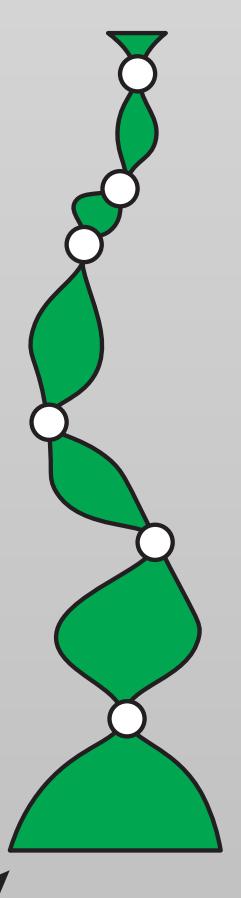
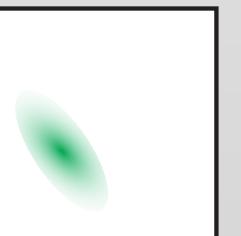
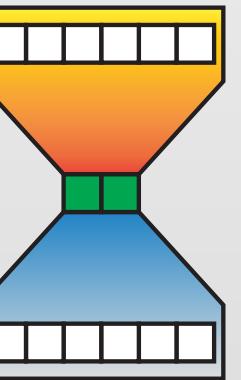
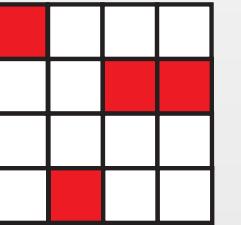
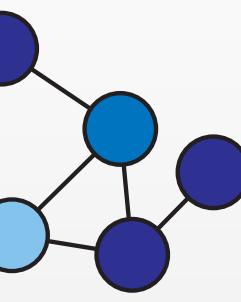


Node Properties

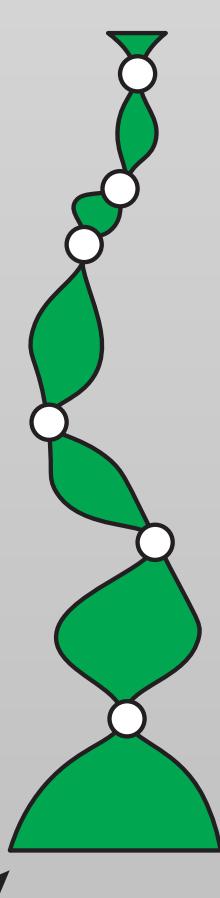
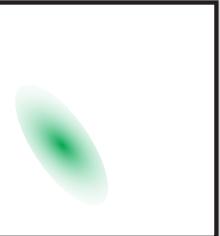
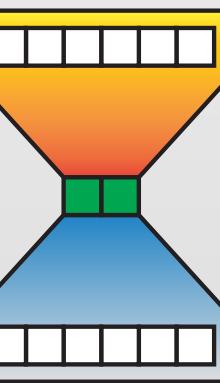
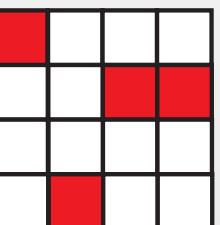
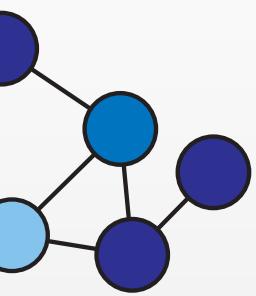
Edge Properties



Forecasting fast-evolving pathogen sequence trajectory

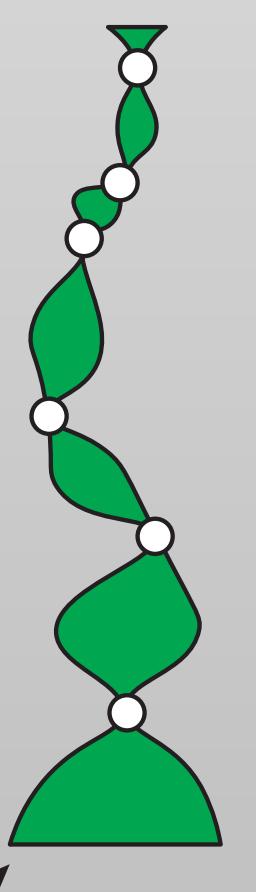
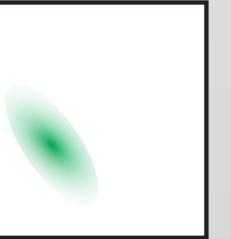
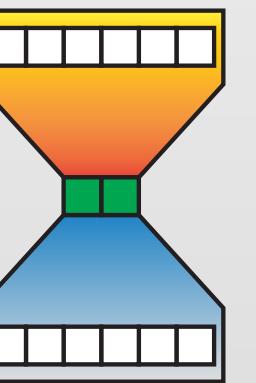
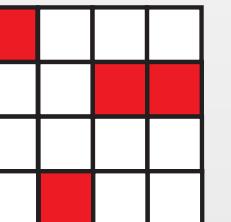
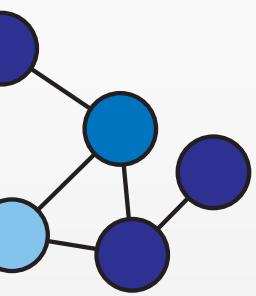


Forecasting sequence evolution is difficult because sequence space is discrete

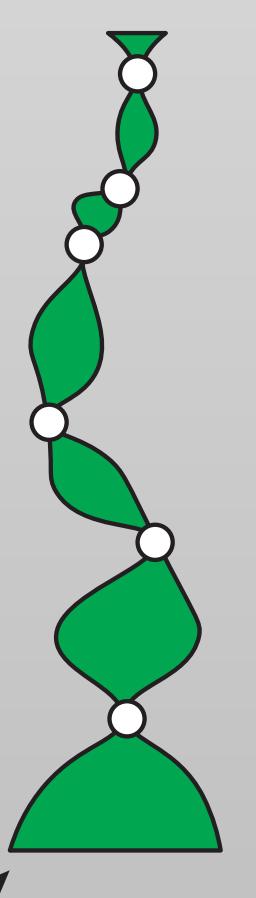
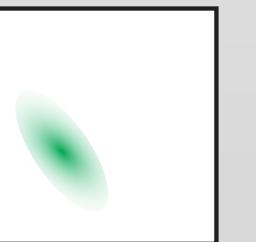
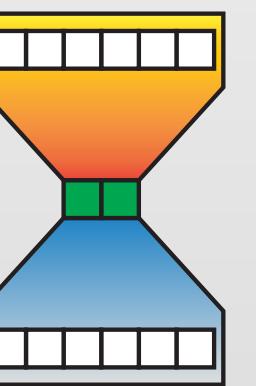
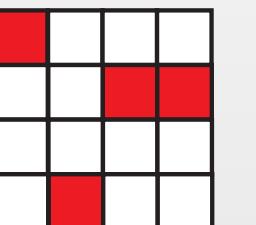
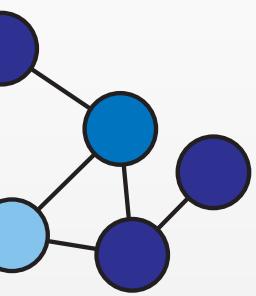


Forecasting sequence evolution is difficult
because sequence space is discrete

Combinatorics problem



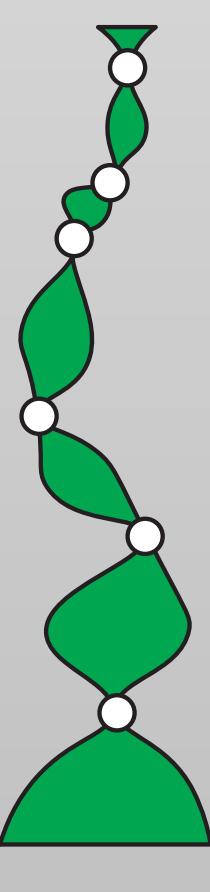
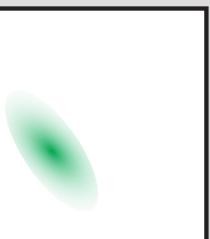
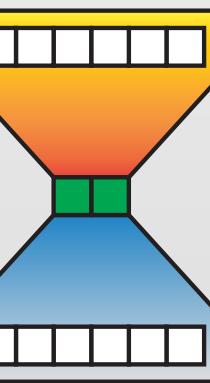
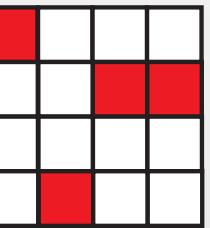
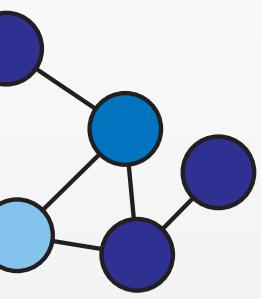
Forecasting sequence evolution is difficult
because sequence space is discrete



Combinatorics problem

10 a.a.: 20^{10} theoretical space

Forecasting sequence evolution is difficult because sequence space is discrete

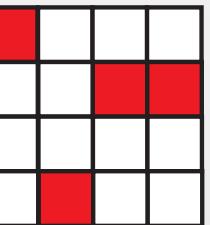
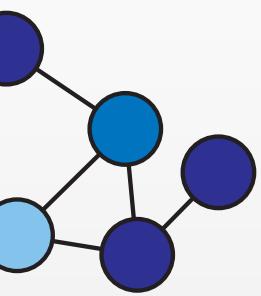


Combinatorics problem

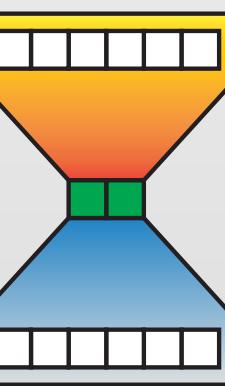
10 a.a.: 20^{10} theoretical space

Deep mutational scans: 10^5 empirical space

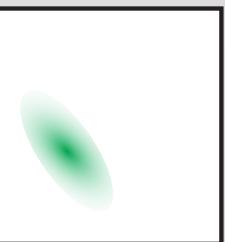
Forecasting sequence evolution is difficult because sequence space is discrete



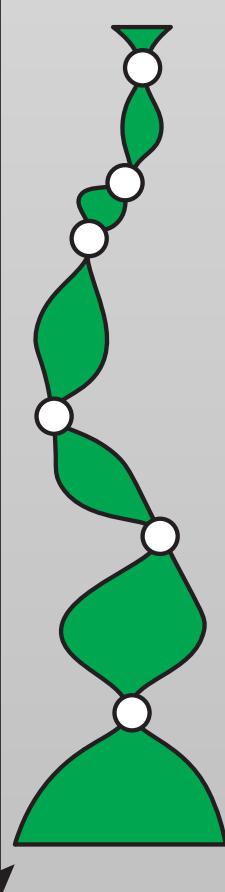
Combinatorics problem



10 a.a.: 20^{10} theoretical space



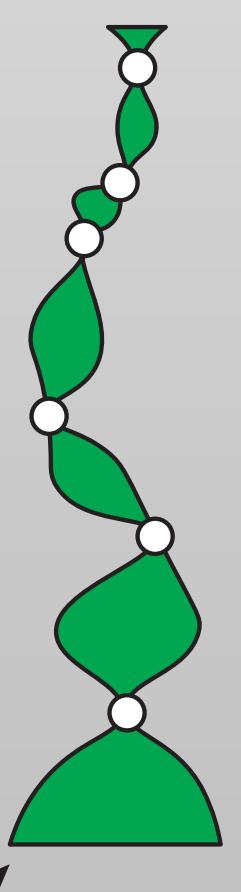
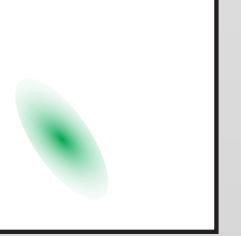
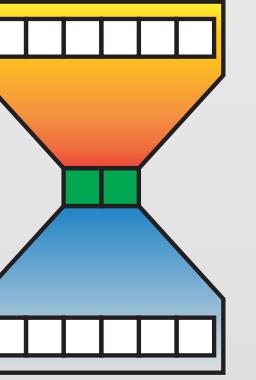
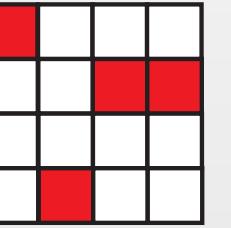
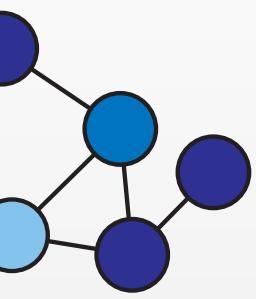
Deep mutational scans: 10^5 empirical space



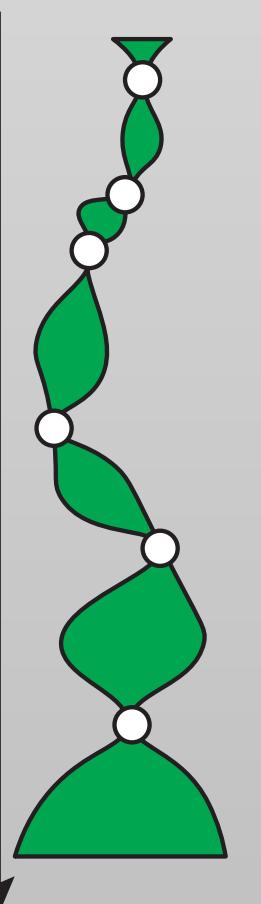
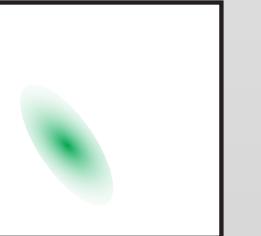
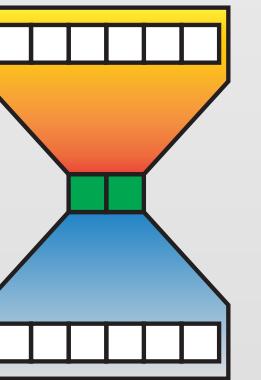
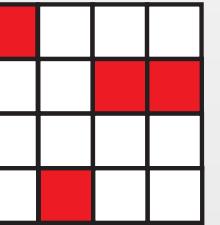
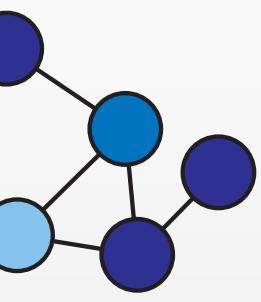
Mutational **direction** cannot be predicted easily

Variational auto-encoders provide a path towards continuous space

SQQ**T**VIPNIGSRPR**VRN**
SQQAVIPNIGSRPRIRD
SQQAVIPNIGSRPRIRD
SQQAVIPNIGSRPRIRD



Variational auto-encoders provide a path towards continuous space



SQQ**T**VIPNIGSRPR**RN**

SQQAVIPNIGSRPRIRD

SQQAVIPNIGSRPRIRD

SQQAVIPNIGSRPRIRD

ACDEF**GHIKLMNPQRSTVWY**

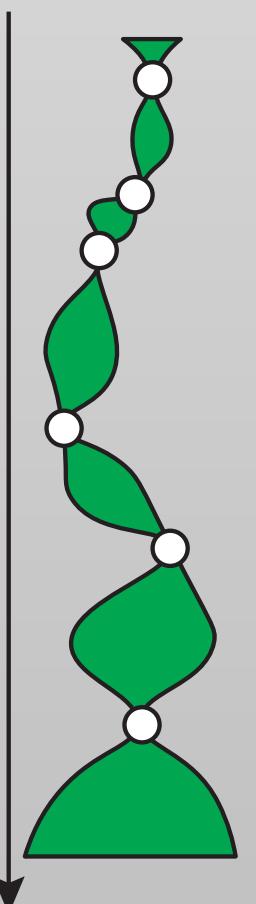
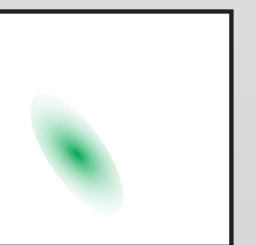
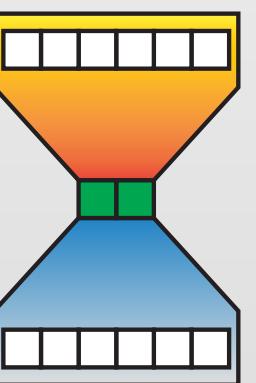
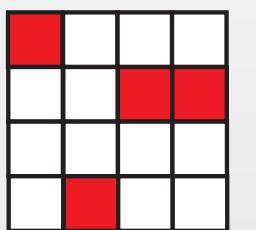
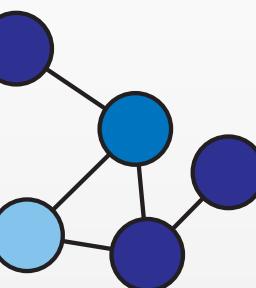
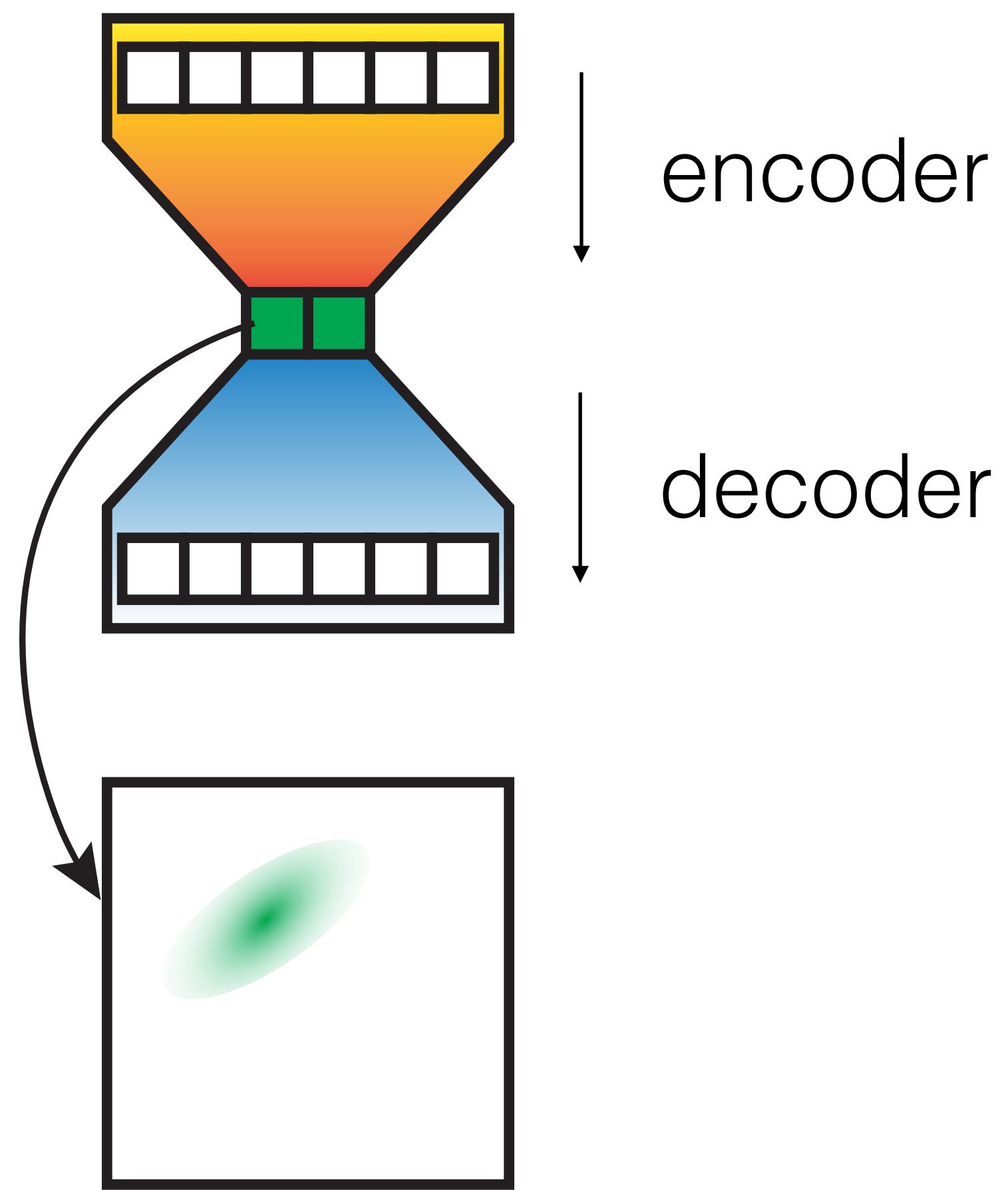
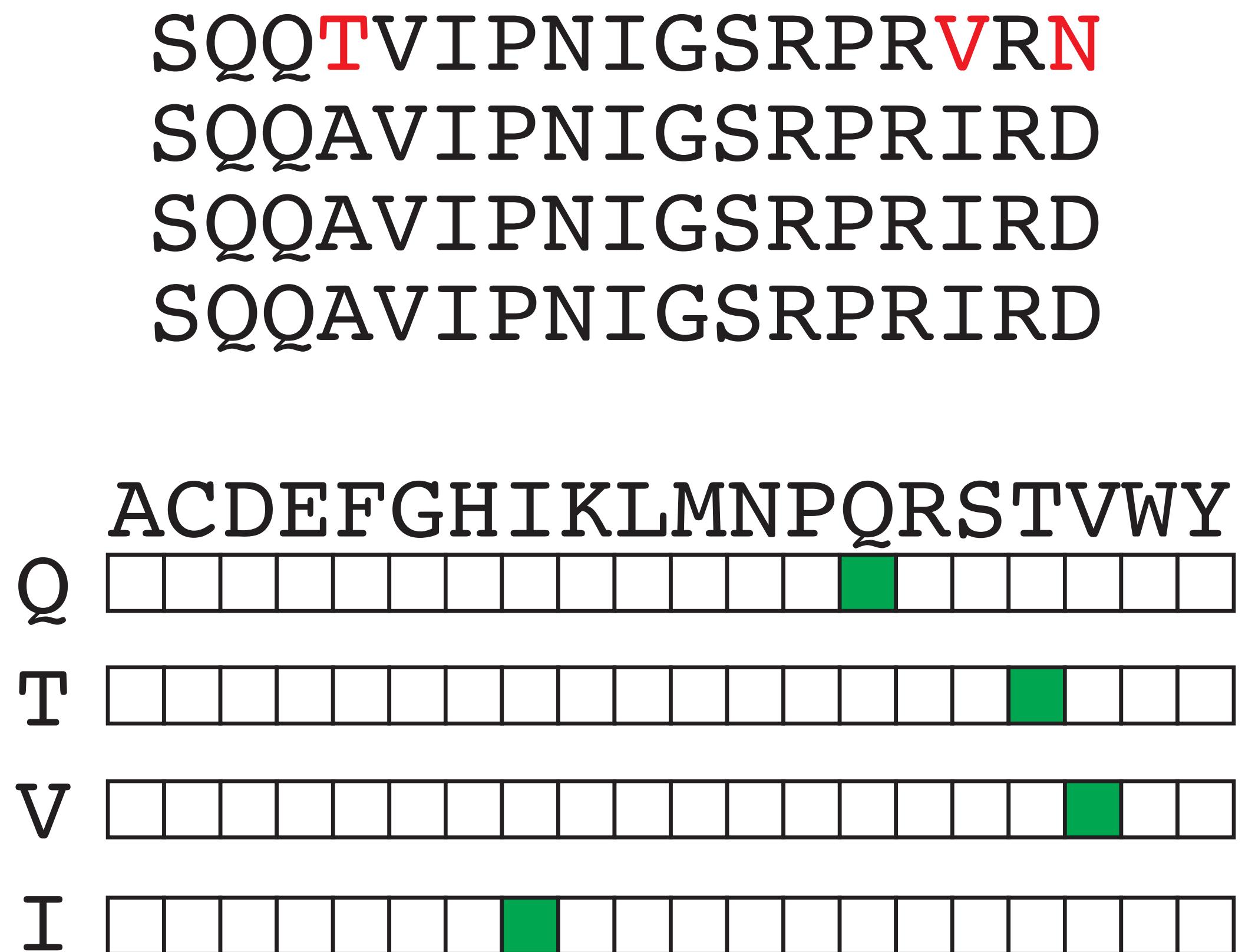
Q A horizontal sequence of 15 boxes. The 13th box from the left is filled green, while the others are white.

T A horizontal sequence of 15 boxes. The 14th box from the left is filled green, while the others are white.

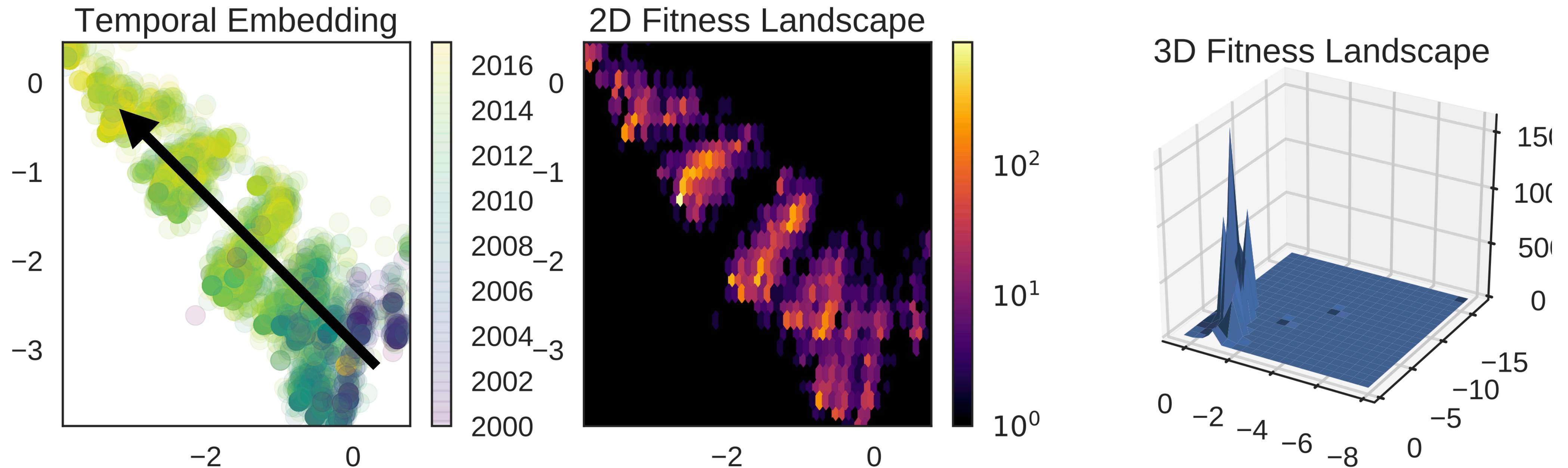
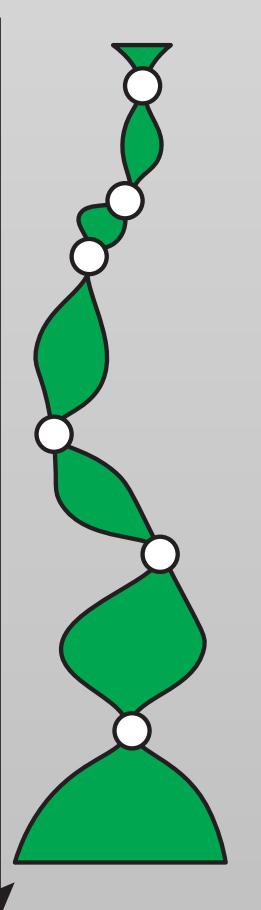
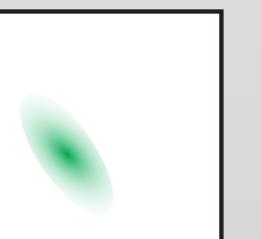
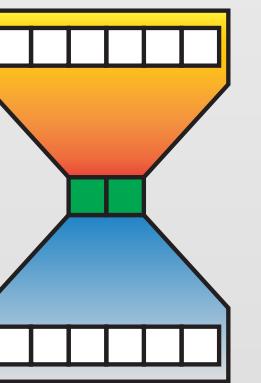
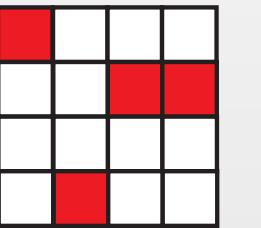
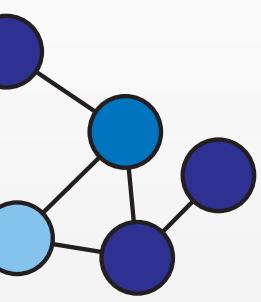
V A horizontal sequence of 15 boxes. The 15th box from the left is filled green, while the others are white.

I A horizontal sequence of 15 boxes. The 12th box from the left is filled green, while the others are white.

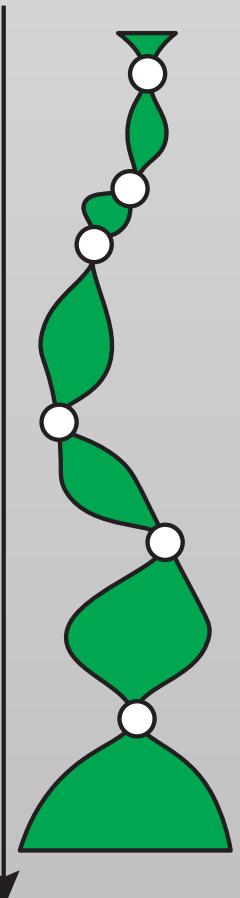
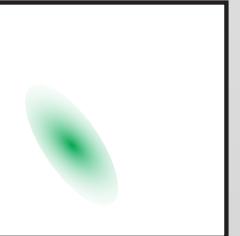
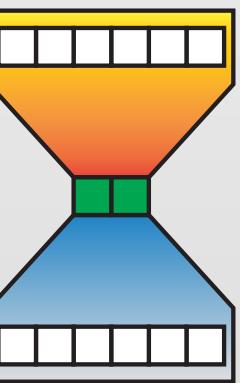
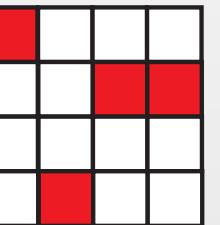
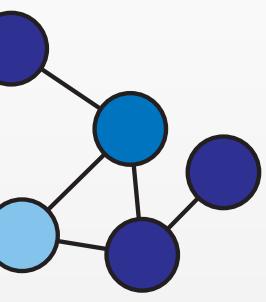
Variational auto-encoders provide
a path towards continuous space



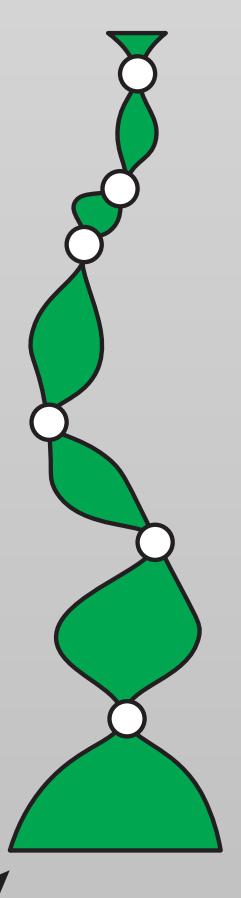
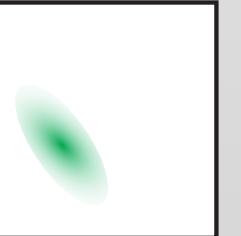
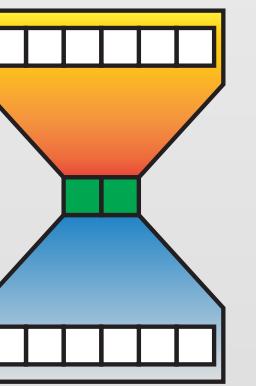
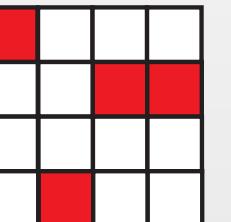
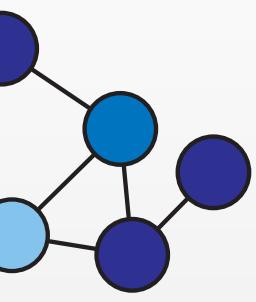
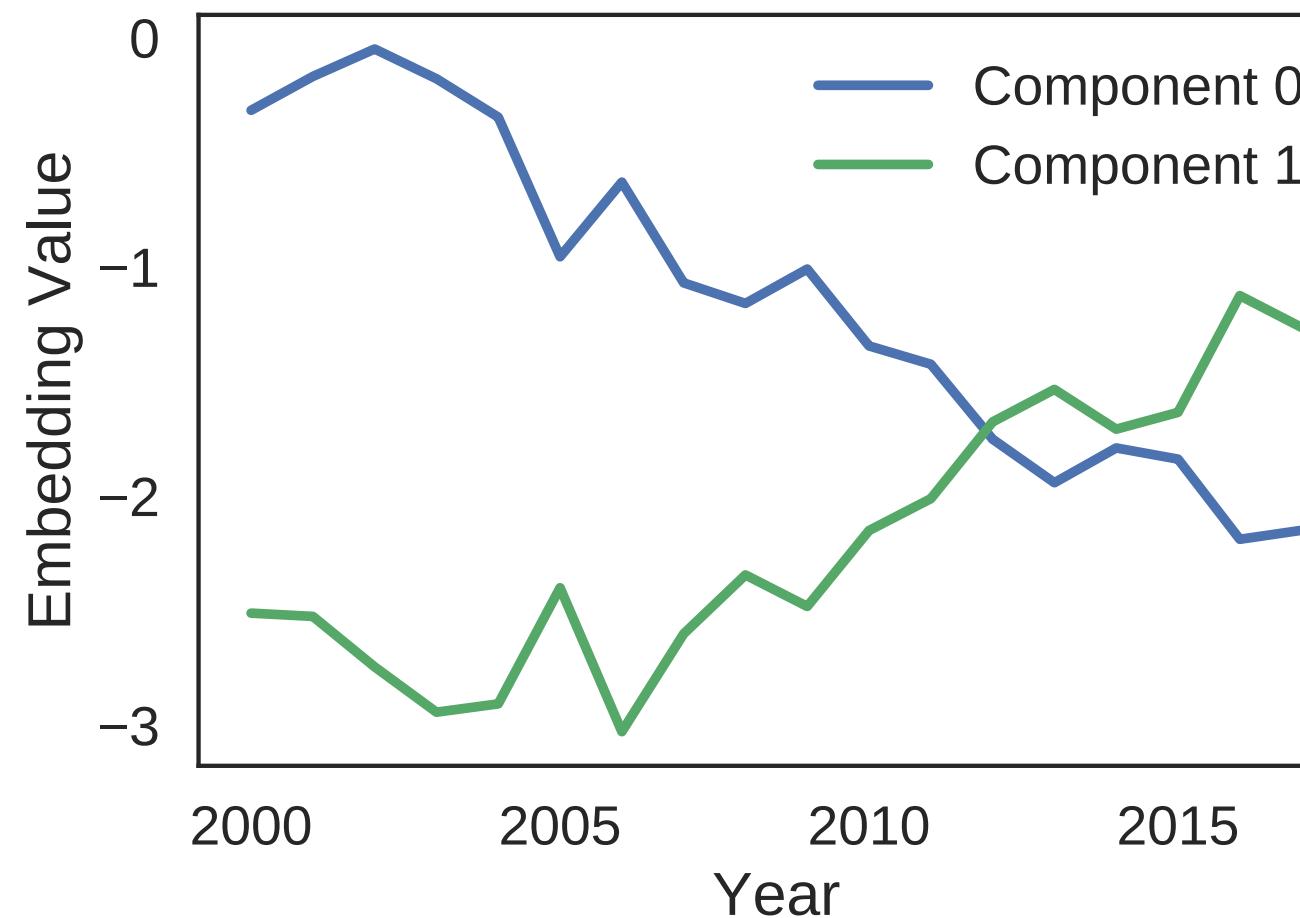
Influenza H3N2 protein evolution has direction in continuous space



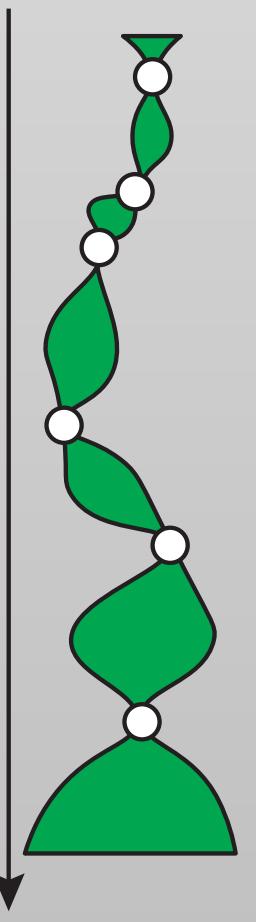
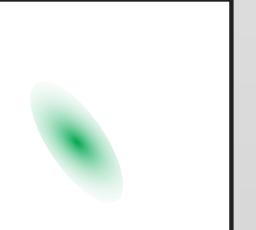
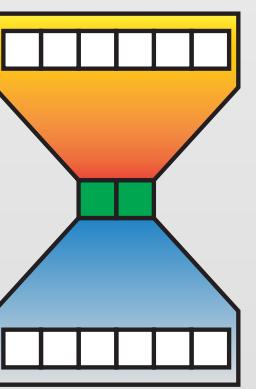
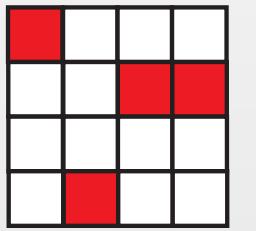
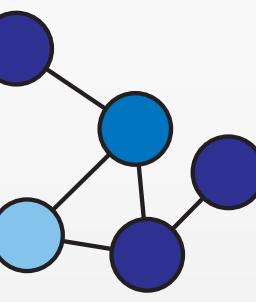
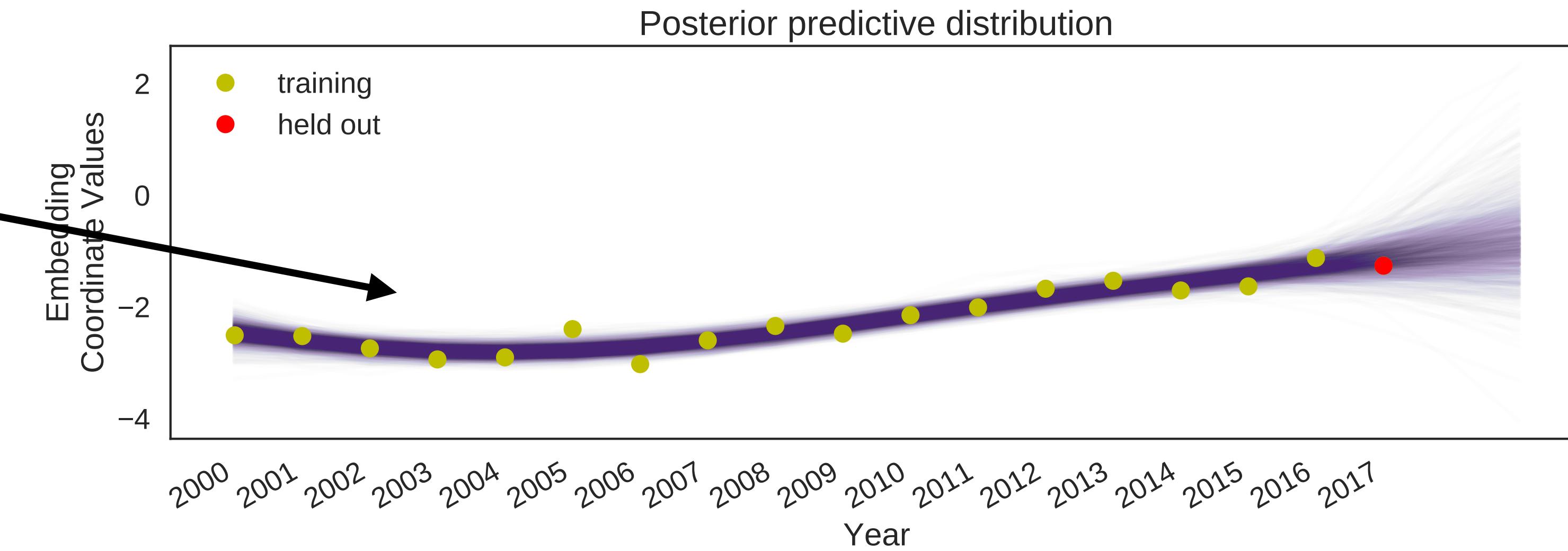
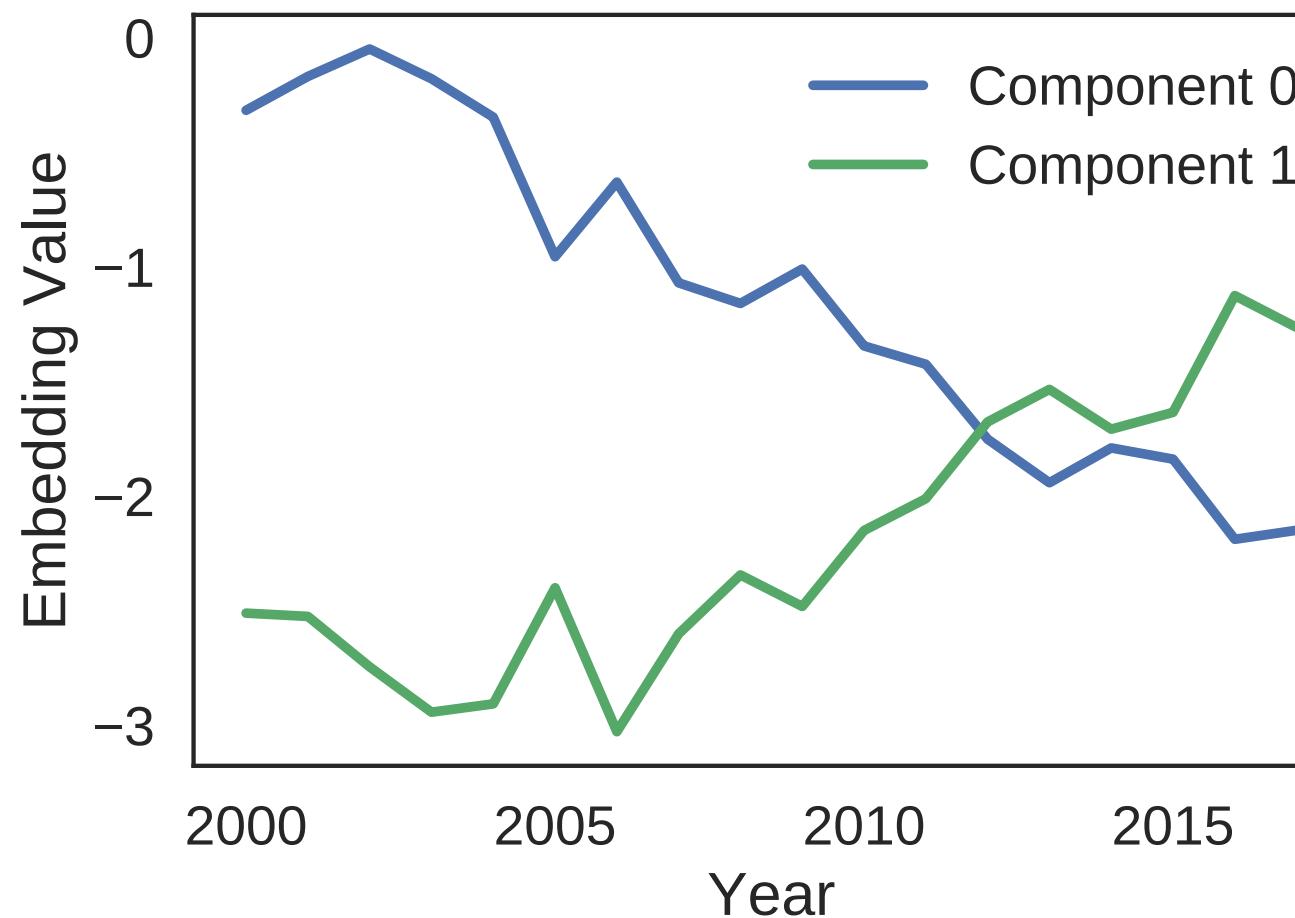
Gaussian process regression enables probabilistic forecasting & sampling of new sequence space



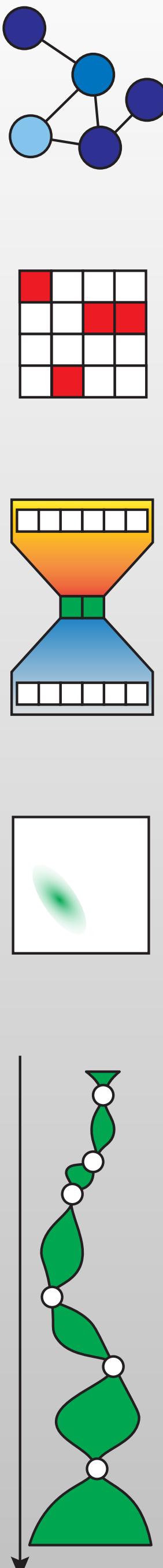
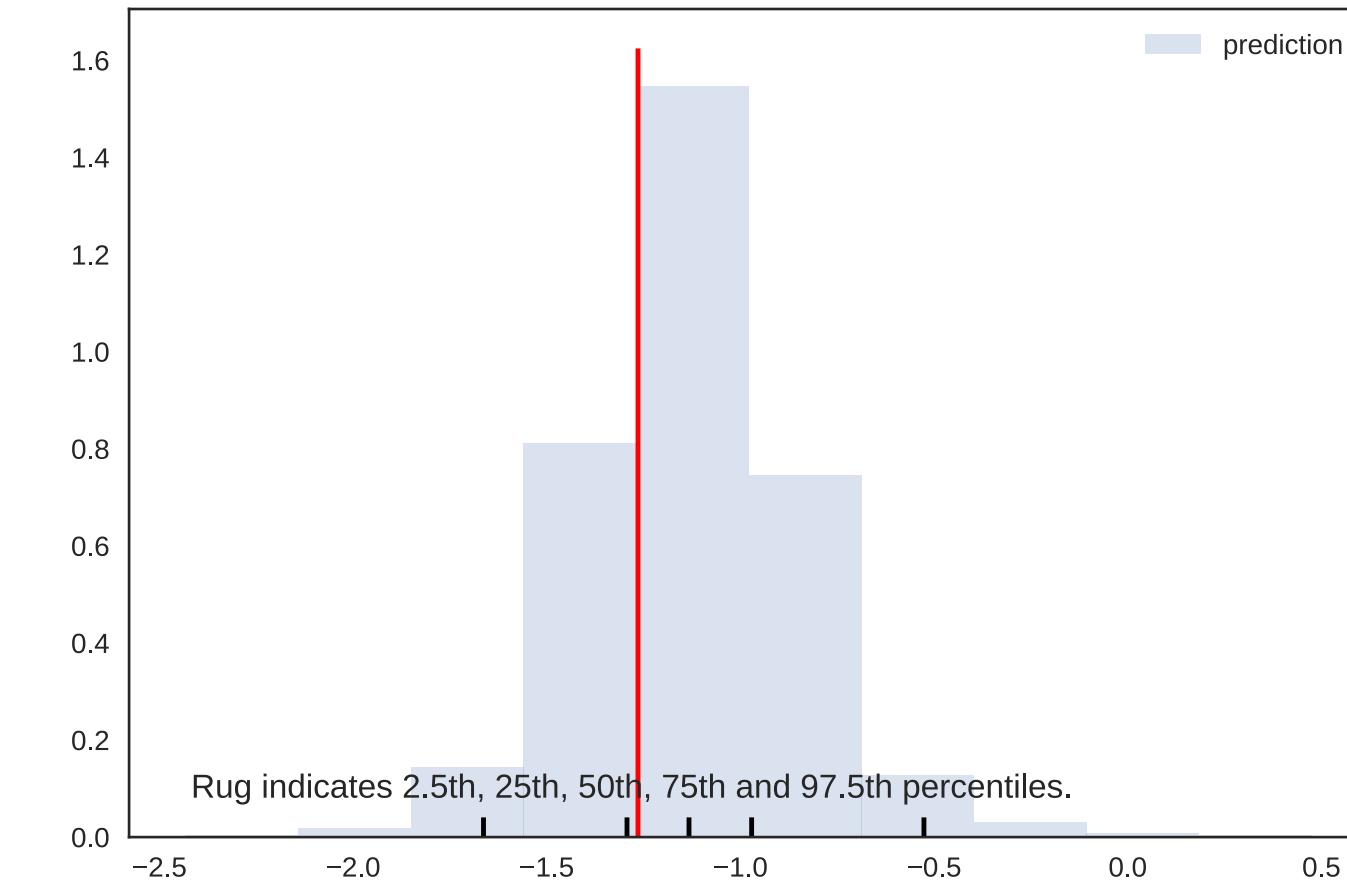
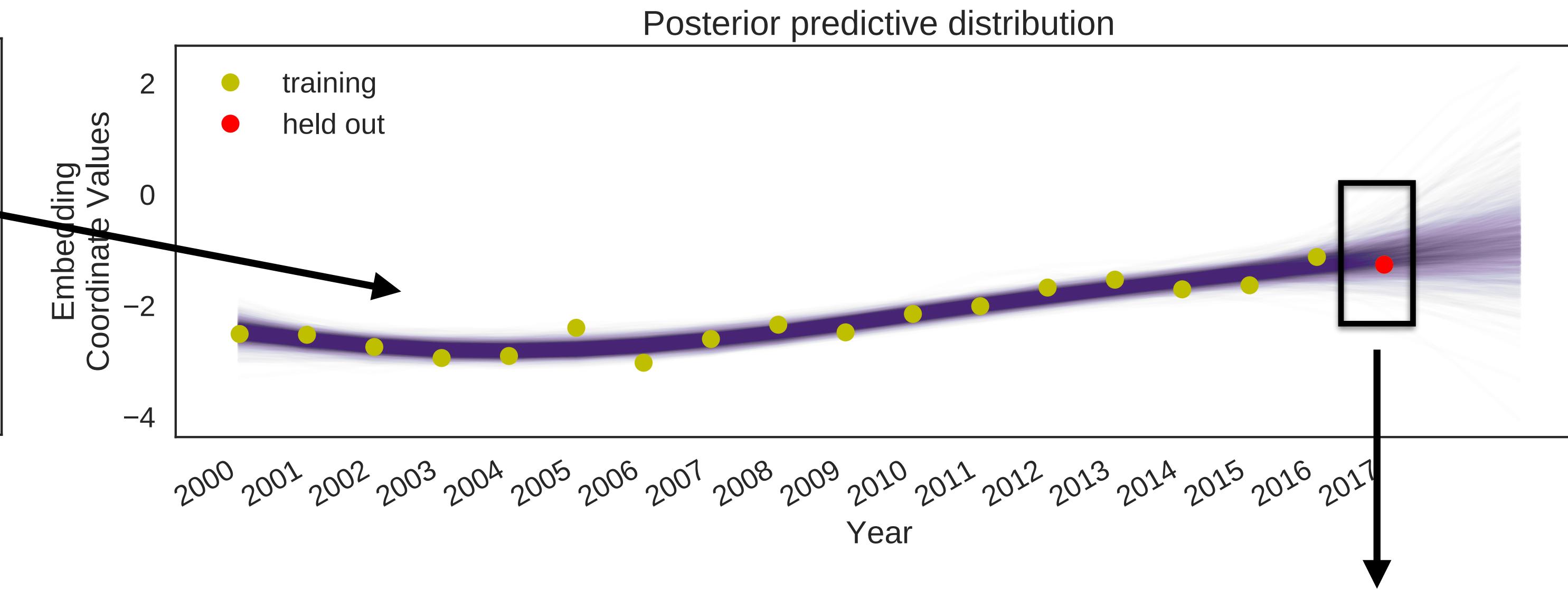
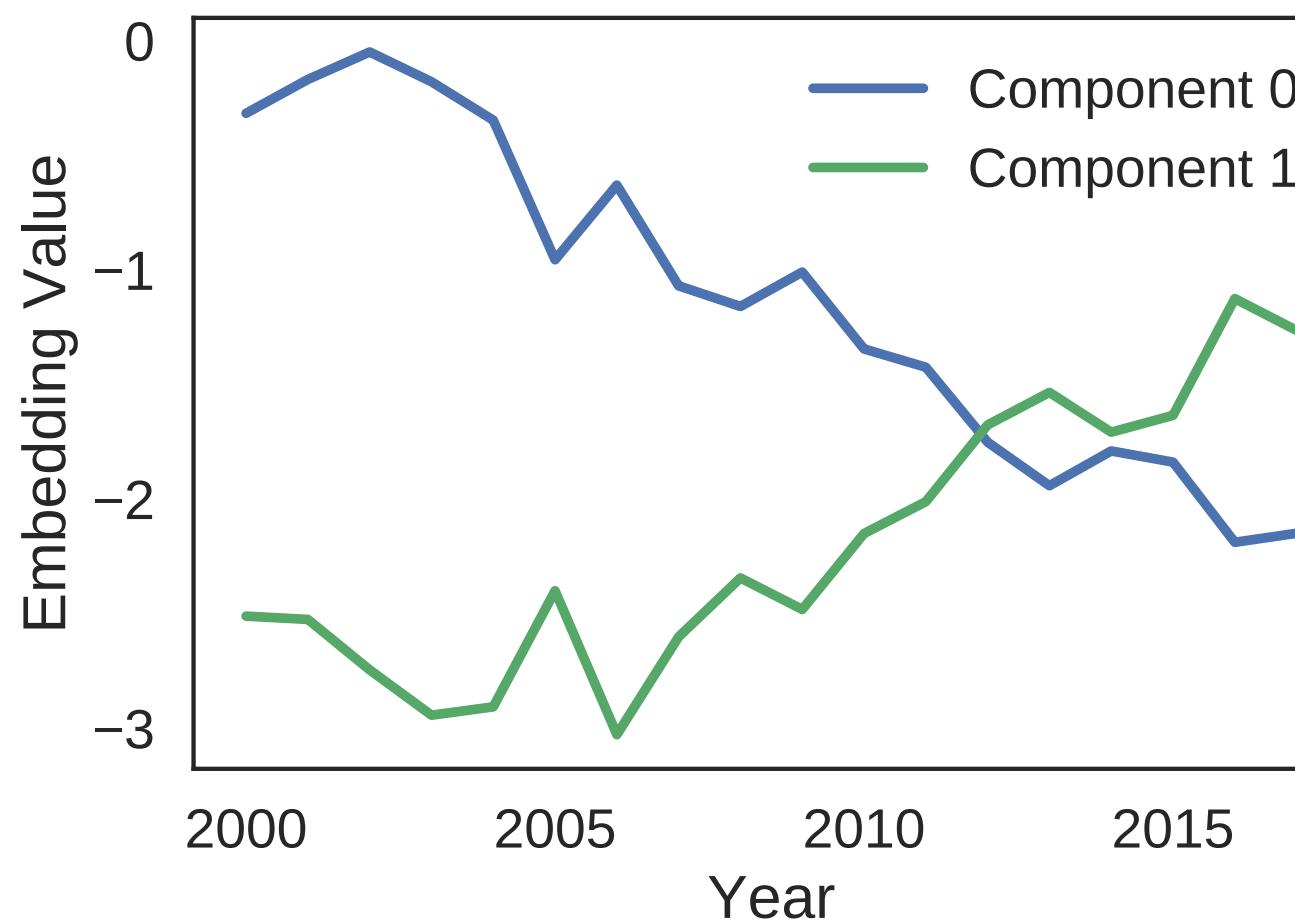
Gaussian process regression enables probabilistic forecasting & sampling of new sequence space



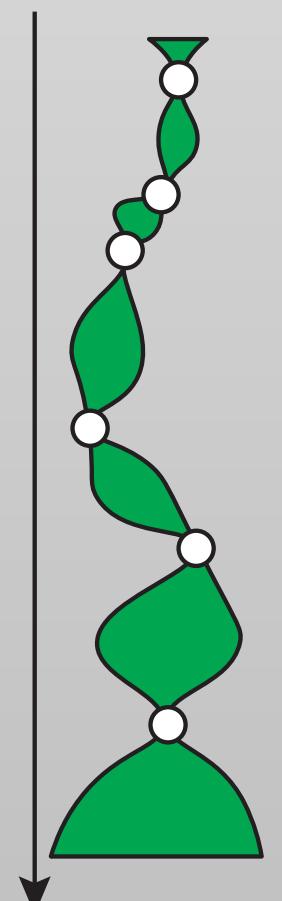
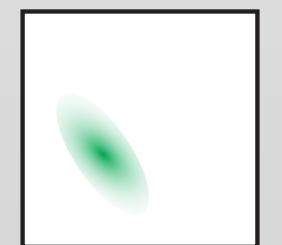
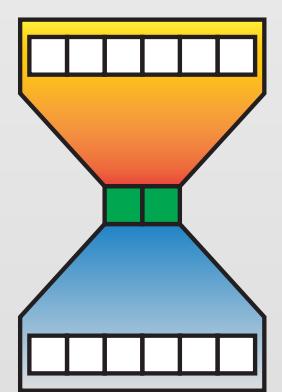
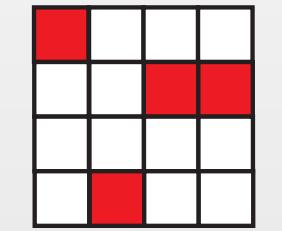
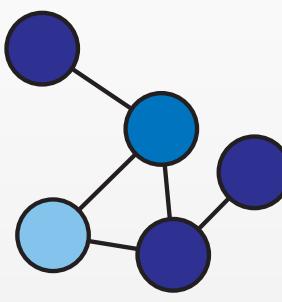
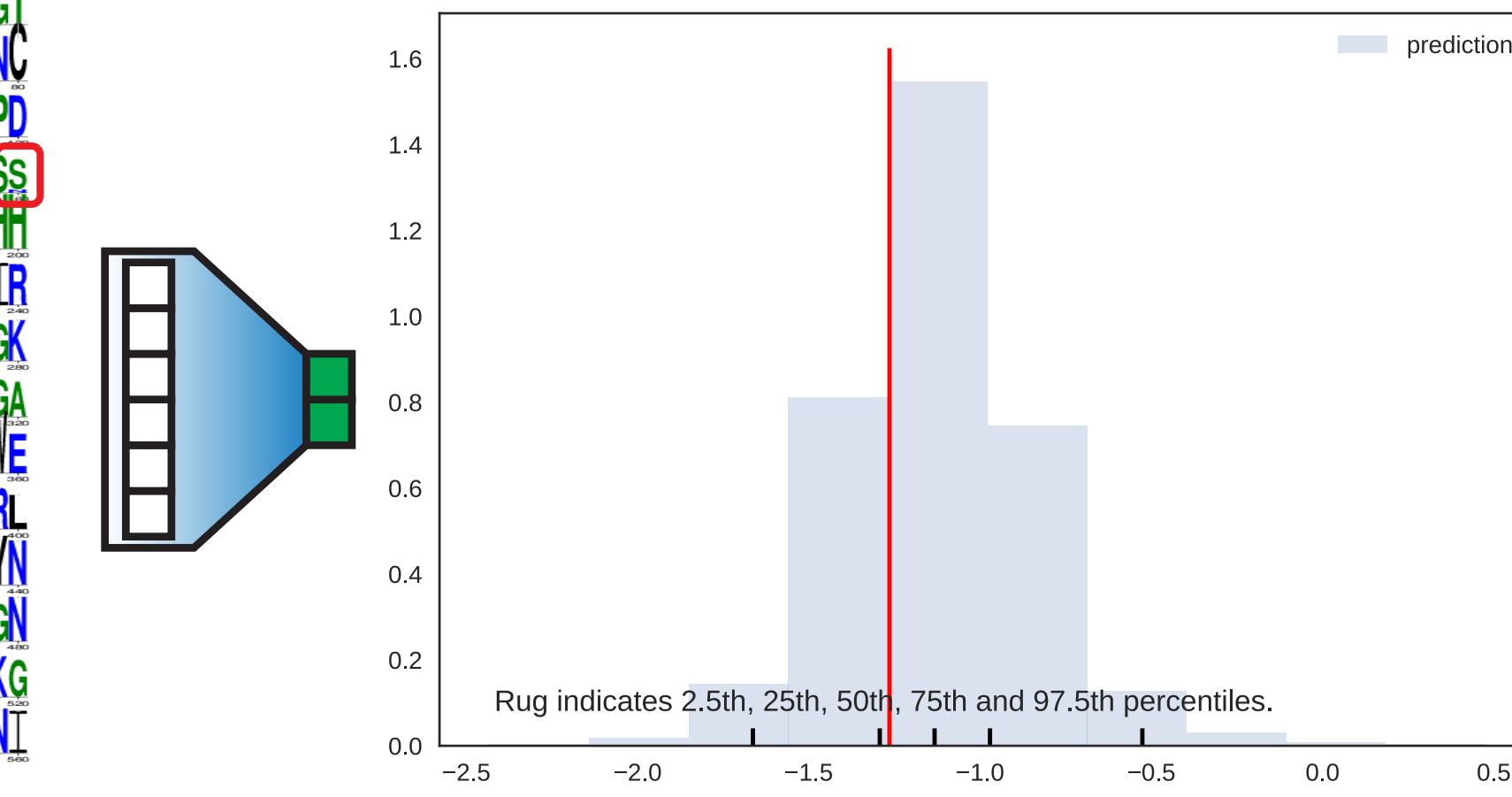
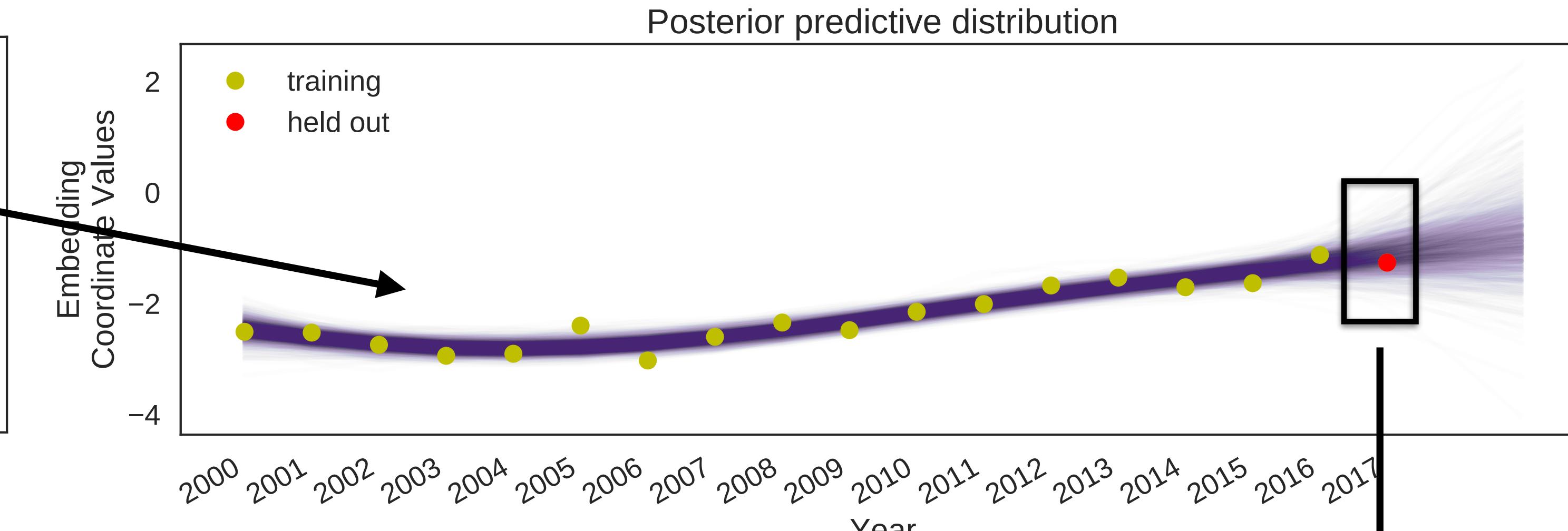
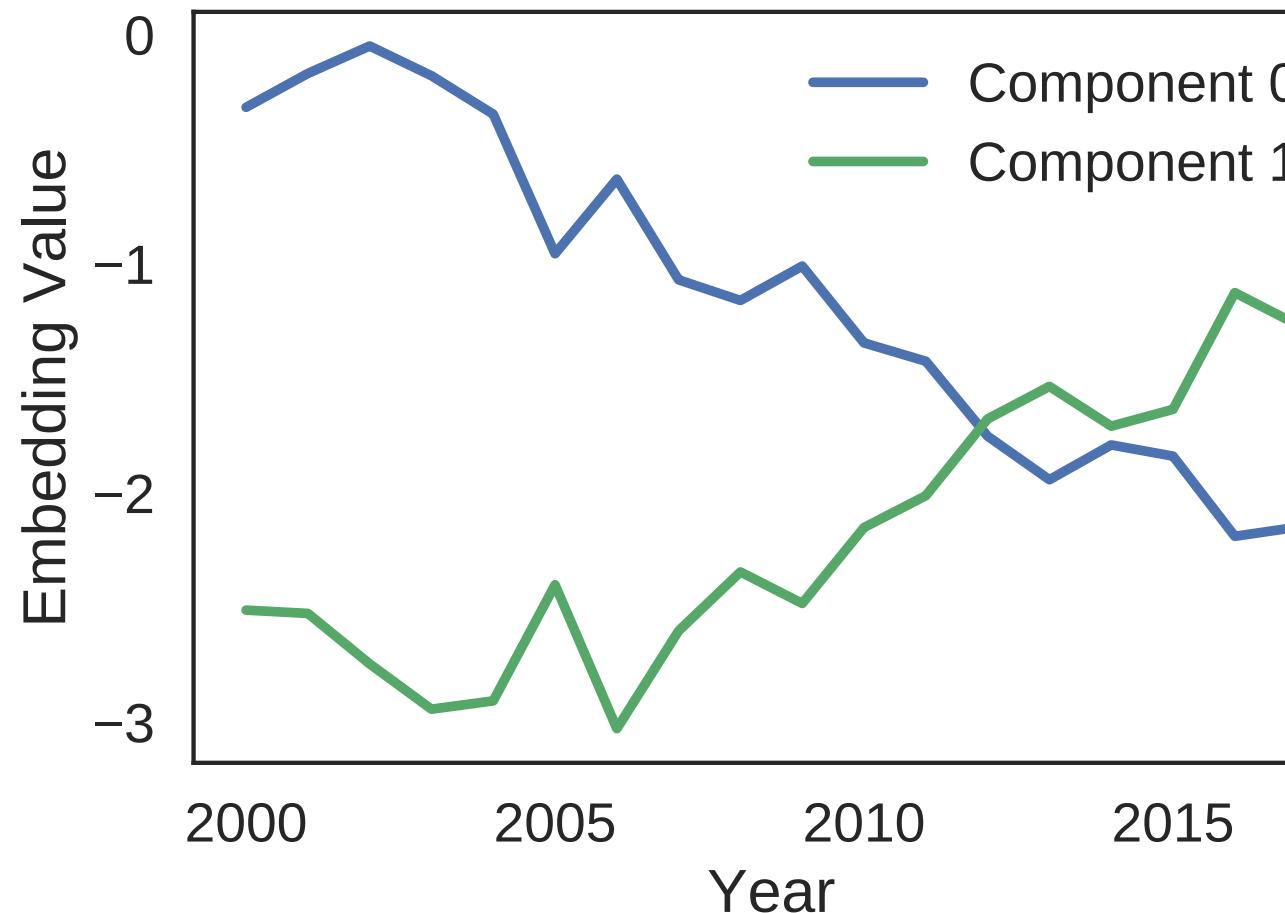
Gaussian process regression enables probabilistic forecasting & sampling of new sequence space



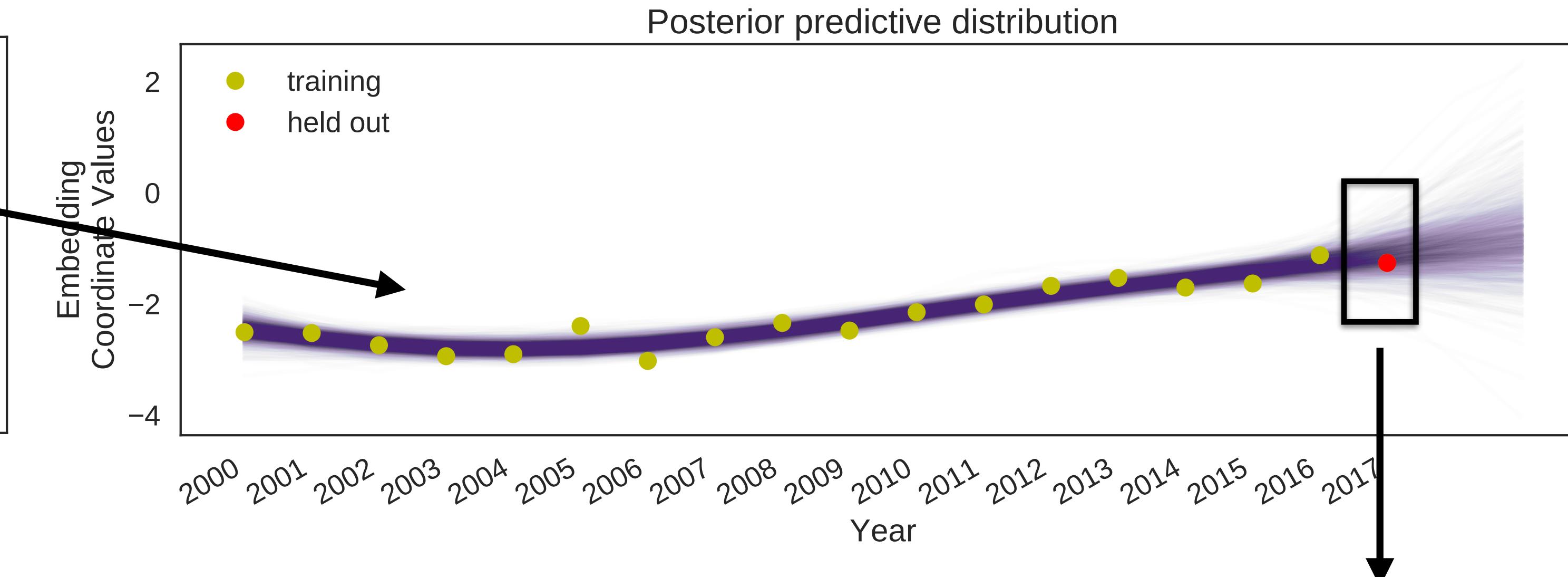
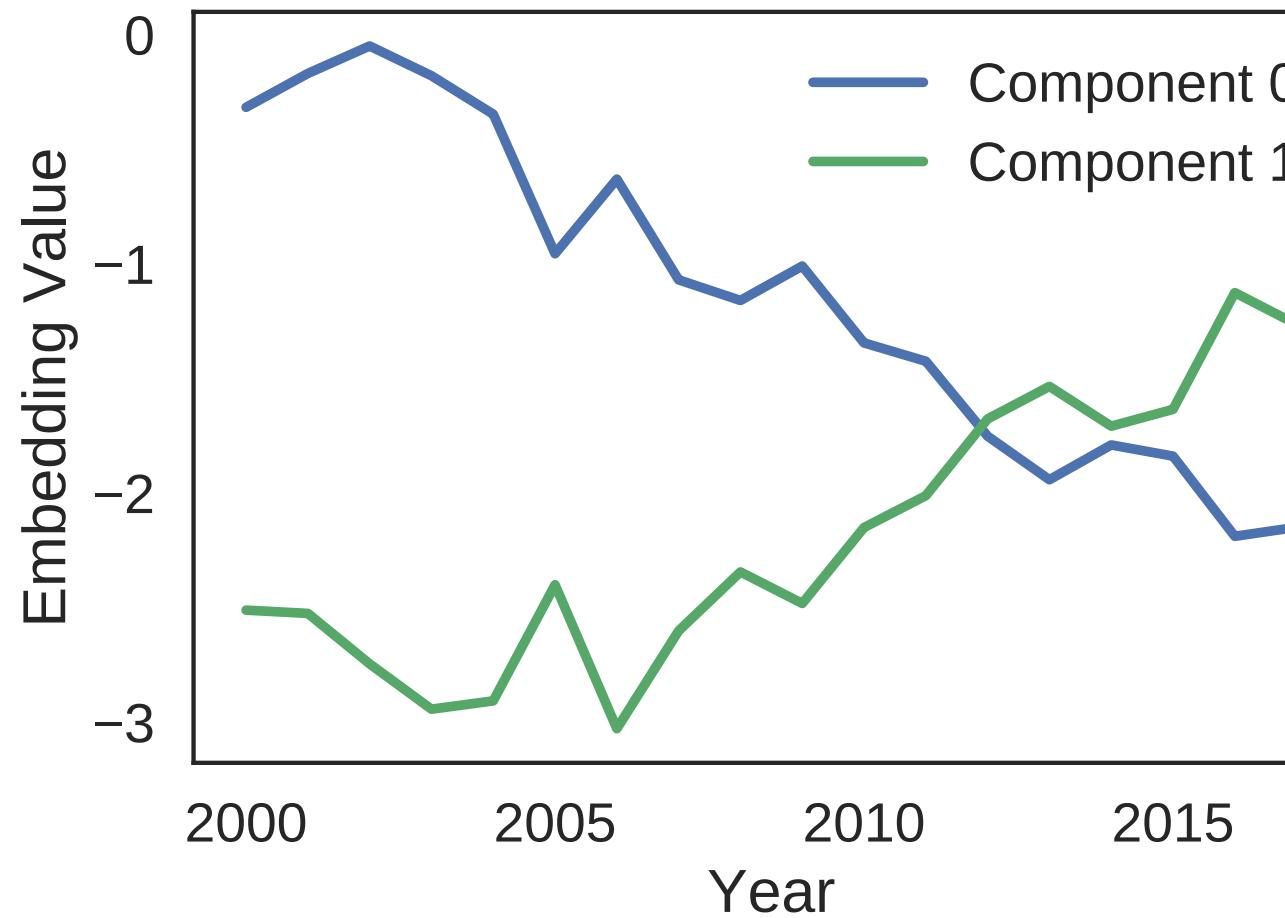
Gaussian process regression enables probabilistic forecasting & sampling of new sequence space



Gaussian process regression enables probabilistic forecasting & sampling of new sequence space

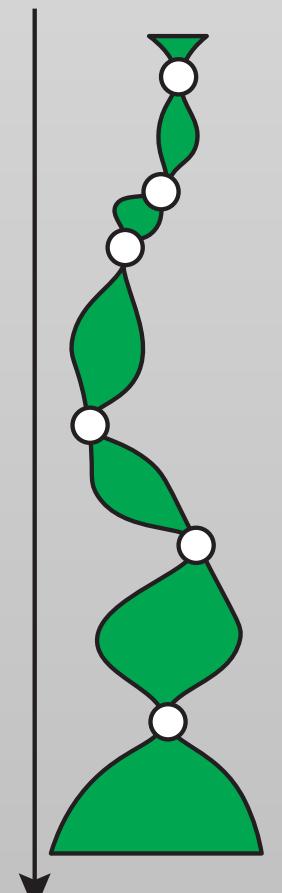
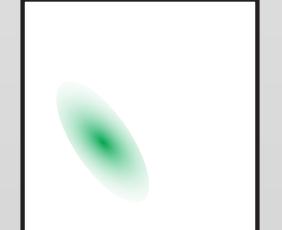
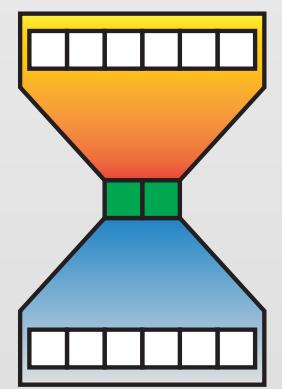
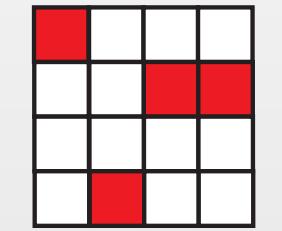
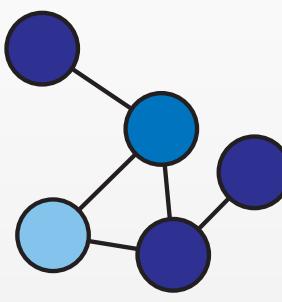
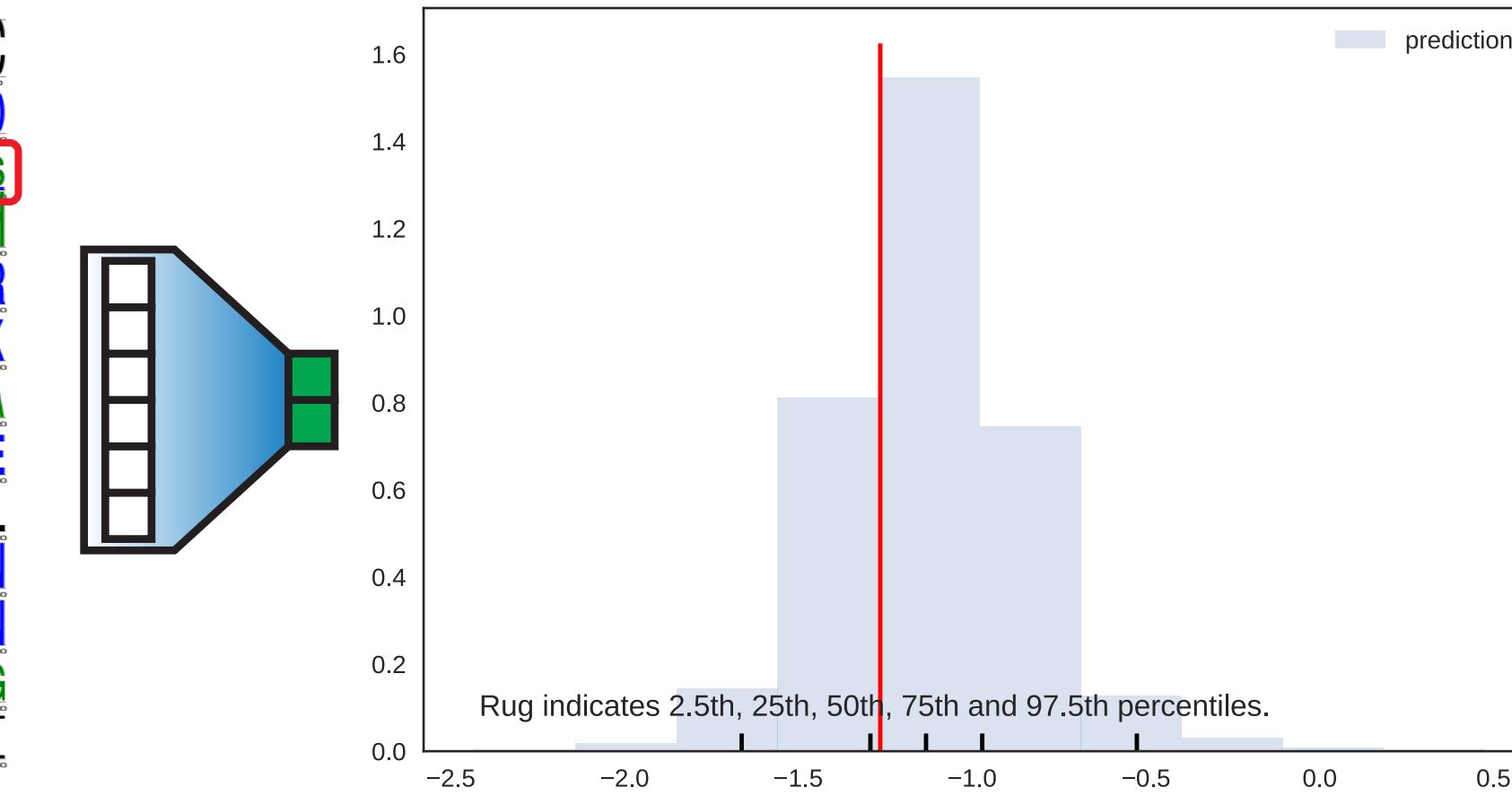
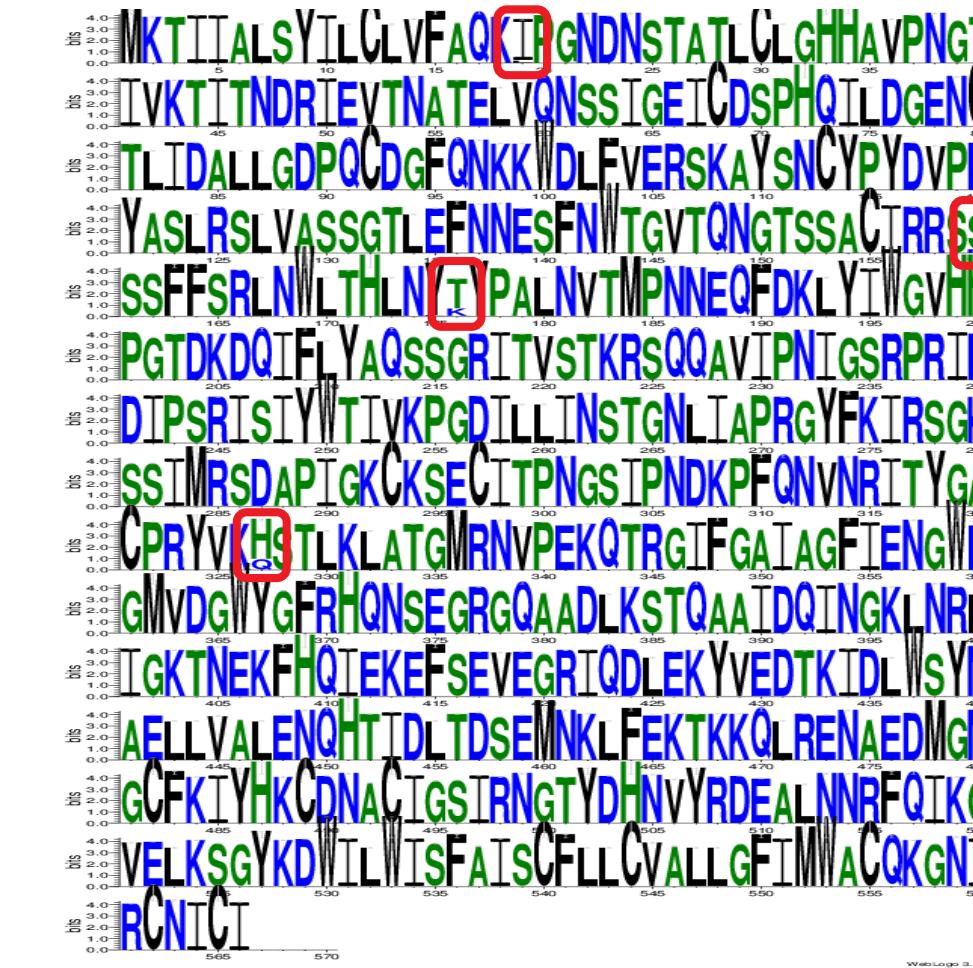


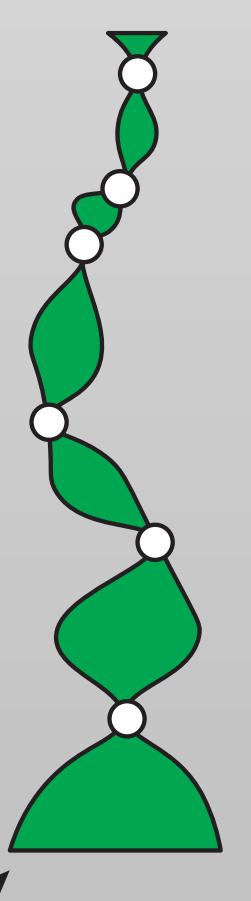
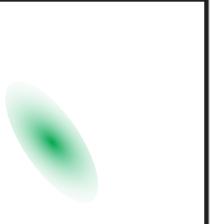
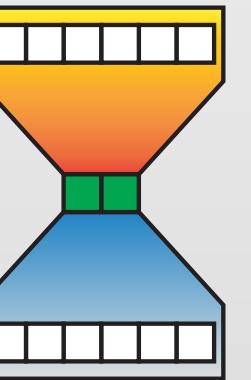
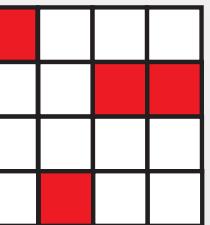
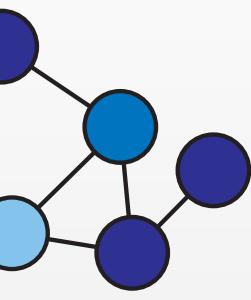
Gaussian process regression enables probabilistic forecasting & sampling of new sequence space



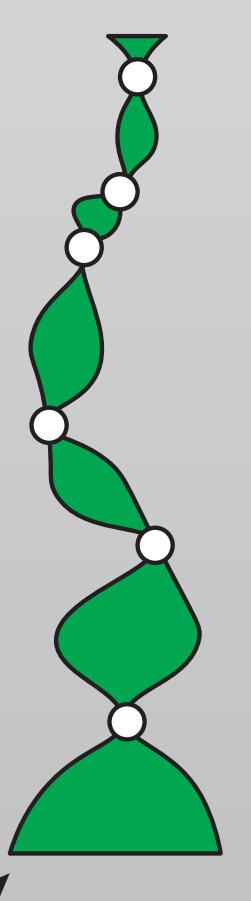
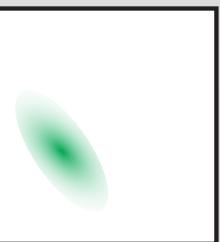
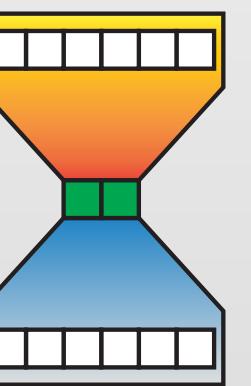
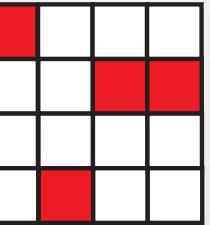
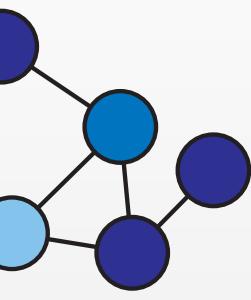
Pos. 176, 327 most likely to change.

Antigenic effects moderate.



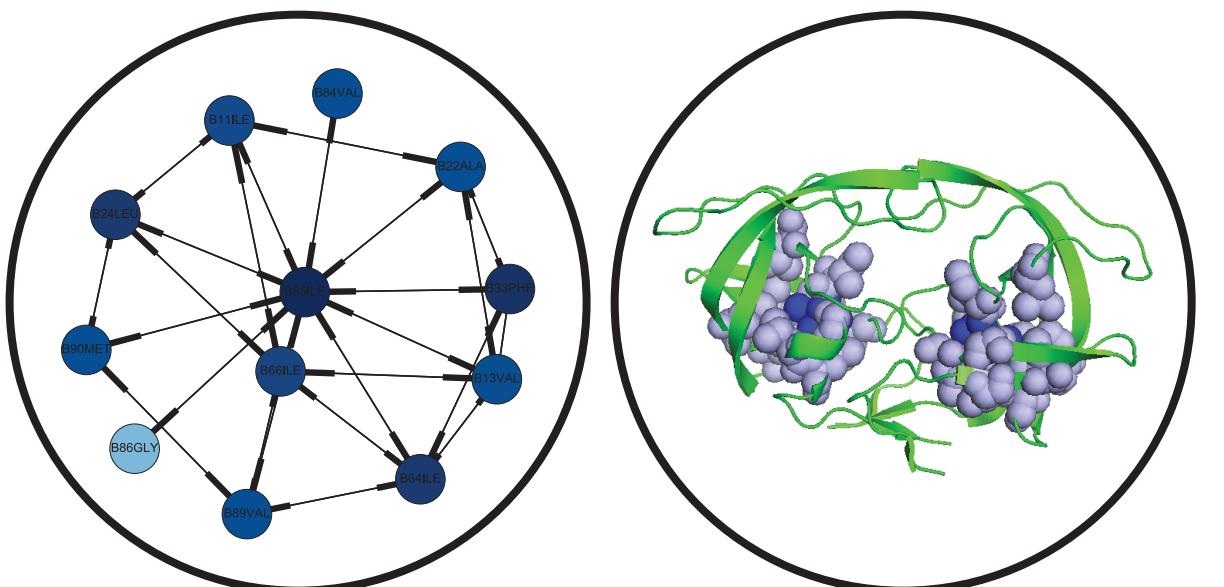


Take-Home Ideas



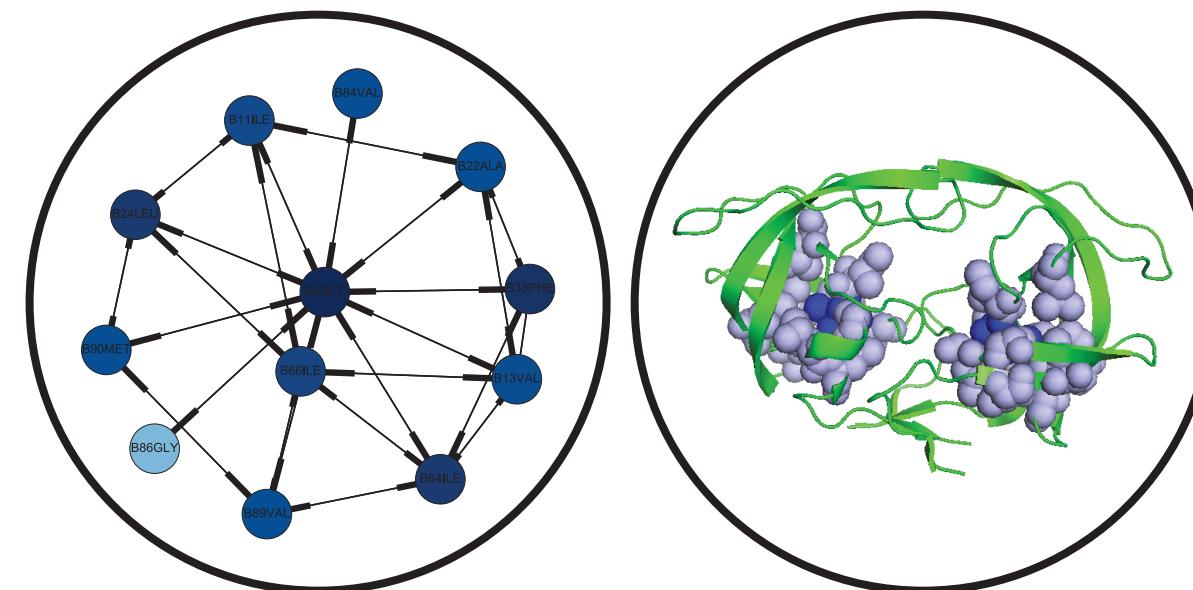
Take-Home Ideas

Graph convolutions: automatically learning structural determinants of phenotype

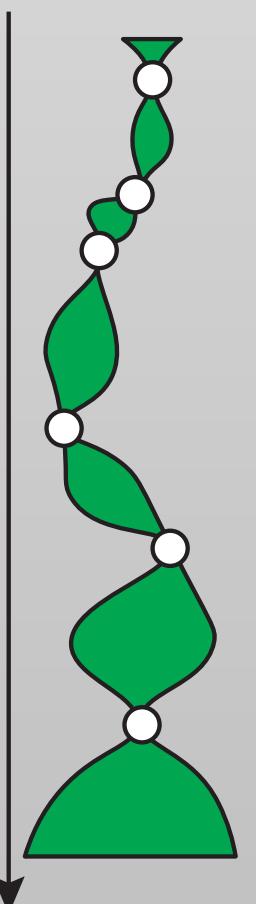
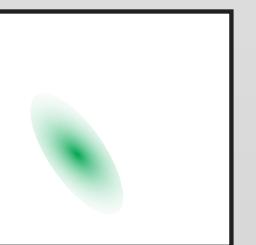
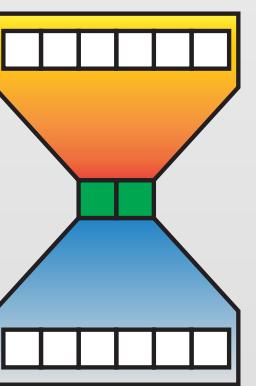
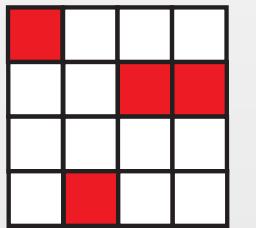
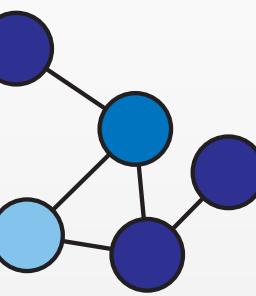
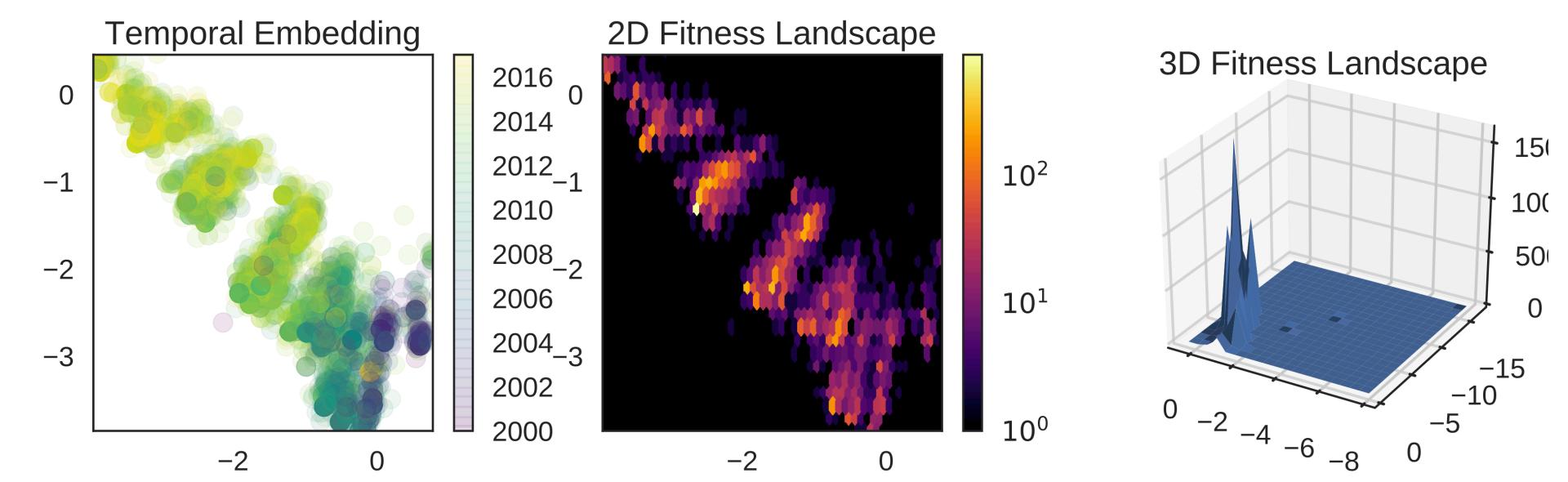


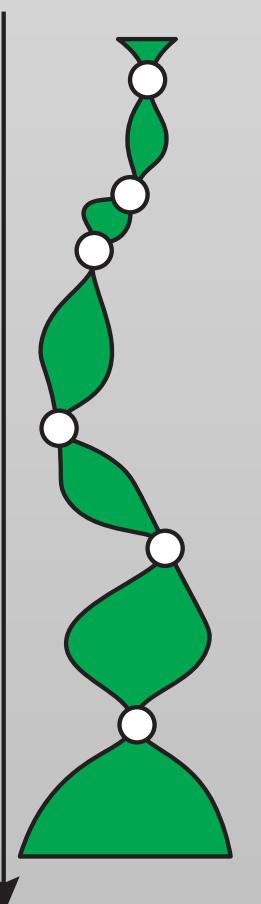
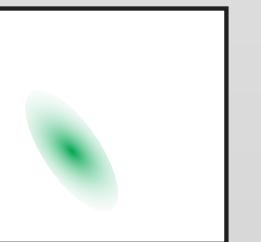
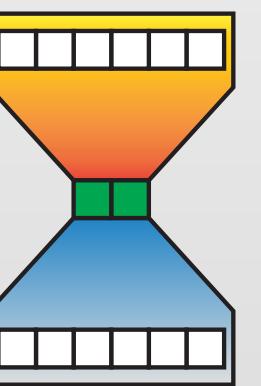
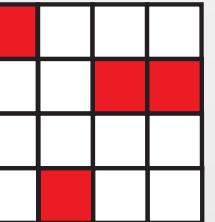
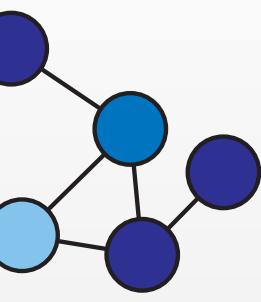
Take-Home Ideas

Graph convolutions: automatically learning structural determinants of phenotype



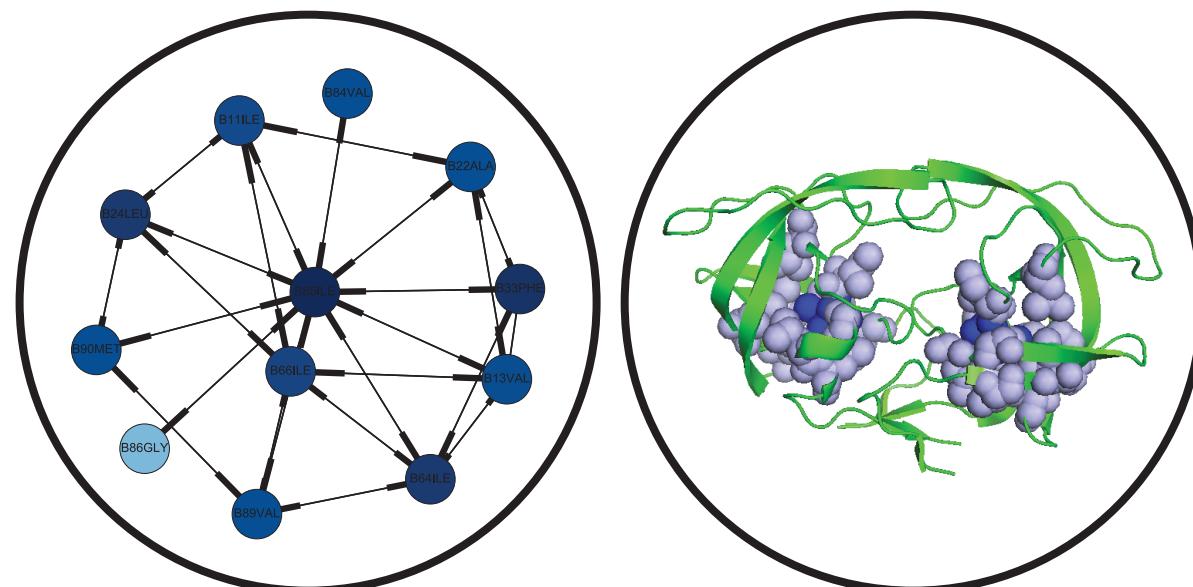
Variational autoencoders: forecasting fast-evolving pathogen sequence trajectory



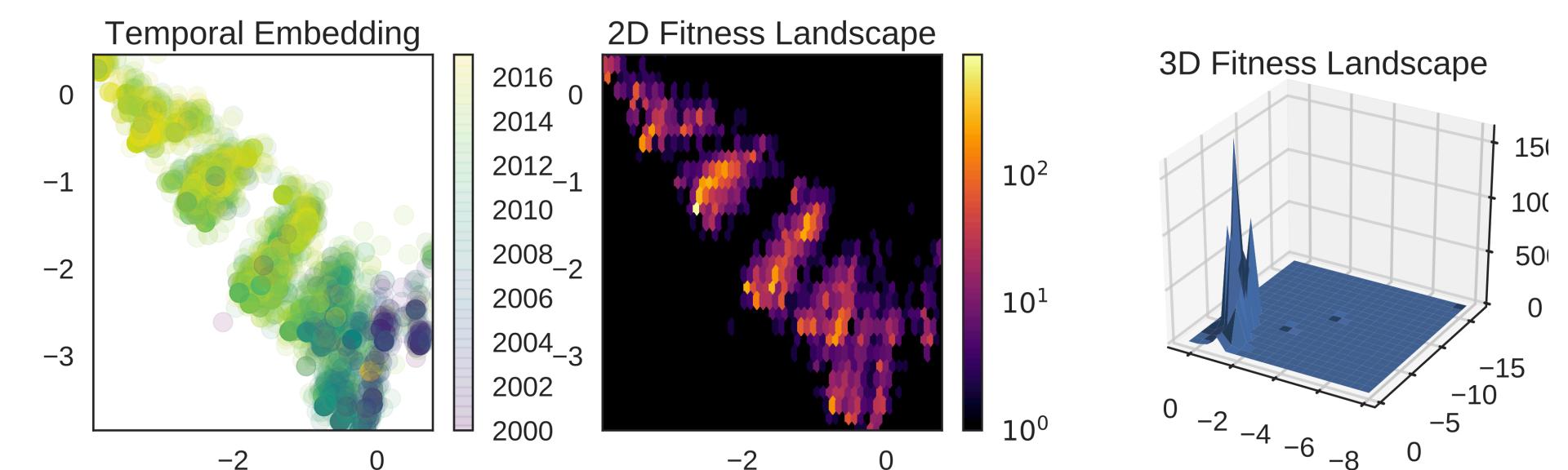


Take-Home Ideas

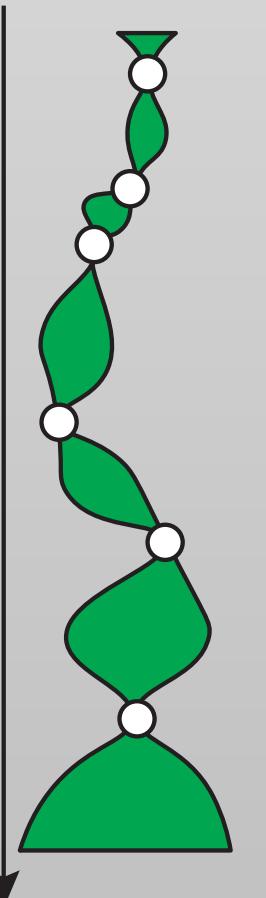
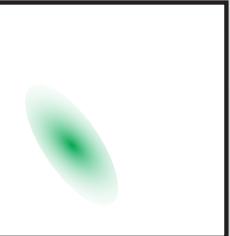
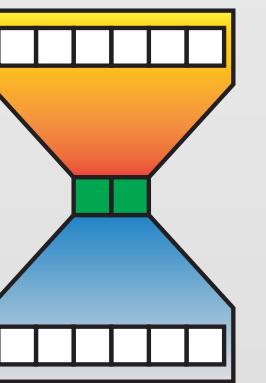
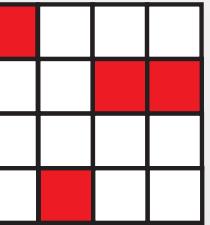
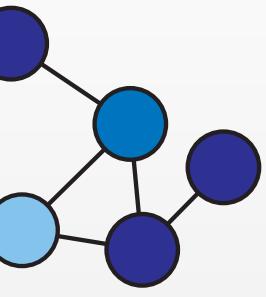
Graph convolutions: automatically learning structural determinants of phenotype



Variational autoencoders: forecasting fast-evolving pathogen sequence trajectory



Applications: CAD of biologics, pre-emptive vaccine development.



Acknowledgments



Jonathan A. Runstadler
Islam T. M. Hussein
Nichola J. Hill

Collaborators
David K. Duvenaud