

Interpretable deep learning models to predict protein phenotype from genotype

Eric J. Ma¹, David K. Duvenaud², Jonathan A. Runstadler^{1,3}

¹Department of Biological Engineering & ³Division of Comparative Medicine, MIT

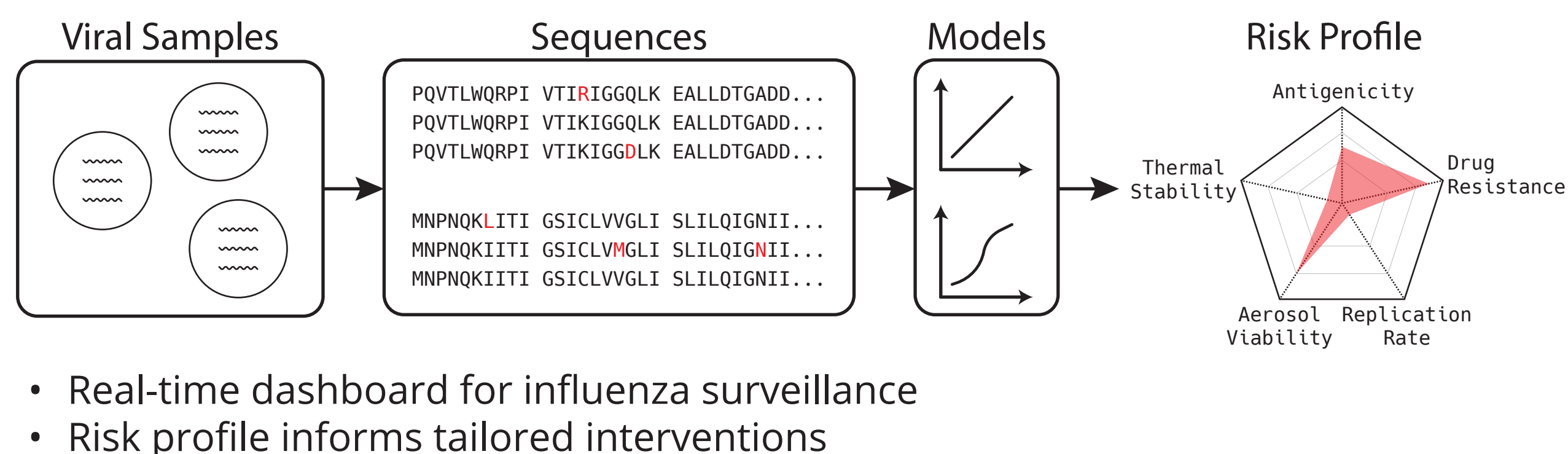
²School of Engineering and Applied Science, Harvard University

Introduction

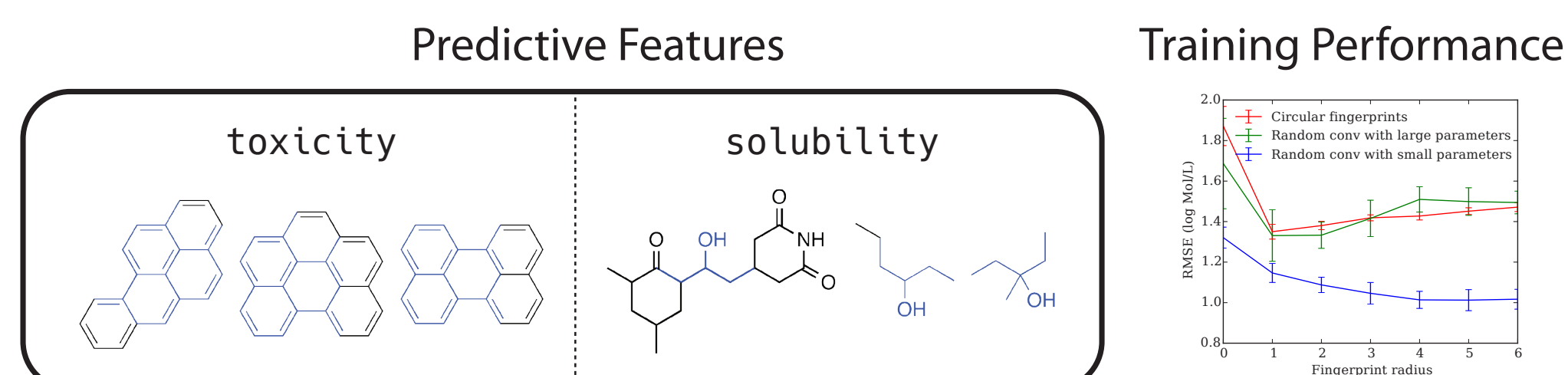
Problems

1. Pathogen risk determination is currently based on non-standardized measurements and simple heuristics.
2. Mapping from genotype to phenotype is complex.
3. Lack of standardized measurements hampers systematic study & reproducibility.
4. Current machine learning models cannot regress on inputs of variable length.

Vision

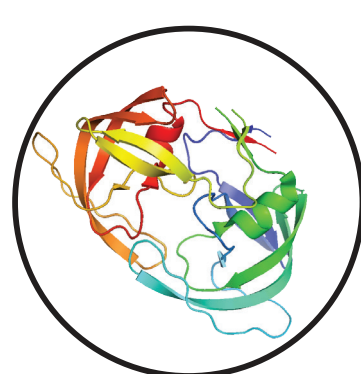


Prior Work



- Duvenaud et. al., 2016 (arXiv): prediction of chemical properties on chemical graphs
- Genotype: chemical structure; phenotype: chemical property.
- Applications in drug screening, toxicity prediction etc.

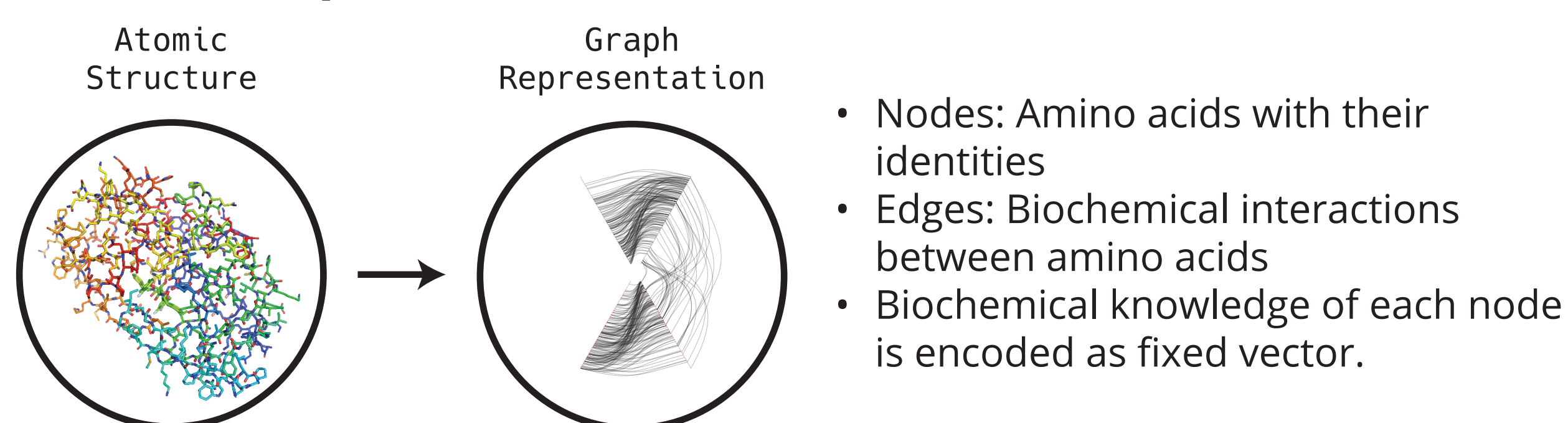
Goals



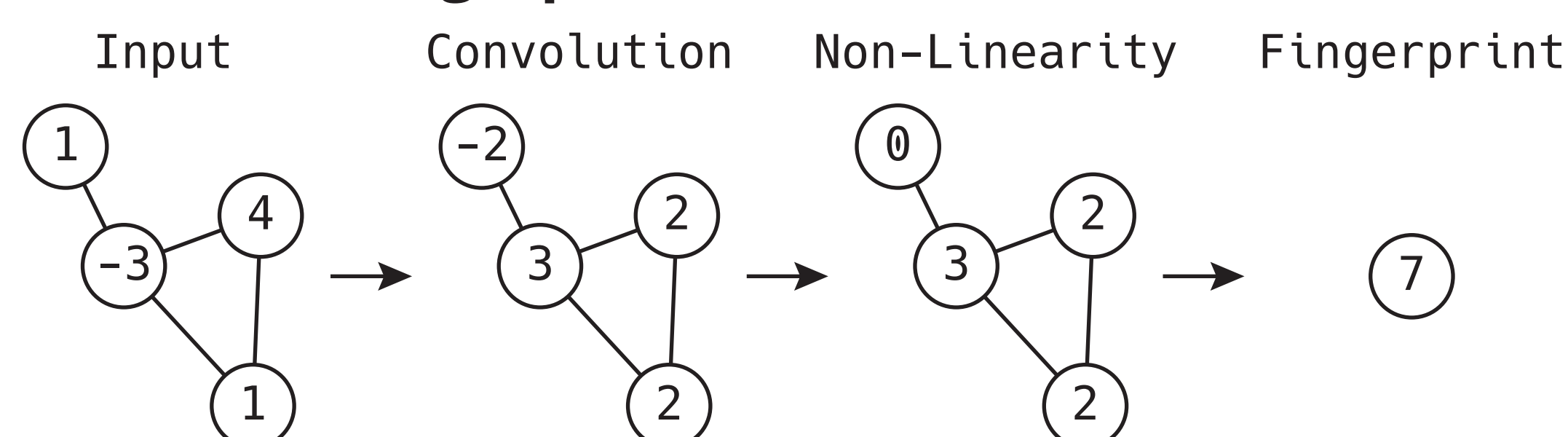
- Data set: HIV-1 Protease, Stanford HIV Drug Resistance Database
- Train convolutional network on protein graph.
- Develop software package for generalized graph regression

Deep Learning Algorithm

Protein Graphs



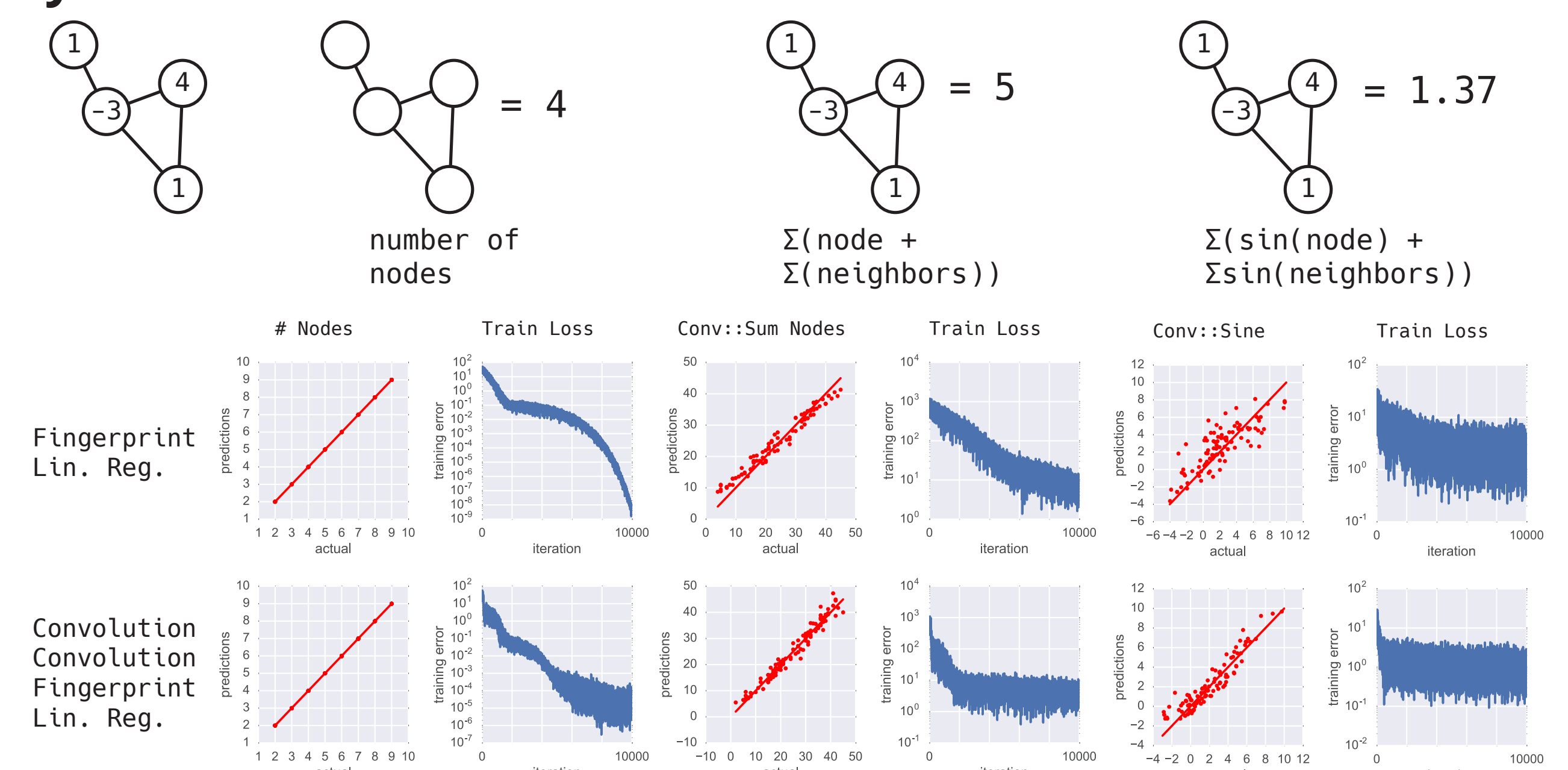
Convolution & Fingerprint



- Convolutions capture local structure of graph
- Non-linearities allow modelling of arbitrary functions
- Fingerprints represent a fixed-length representation of underlying graph.
- Graphs with identical nodes and edges will have identical fingerprints.

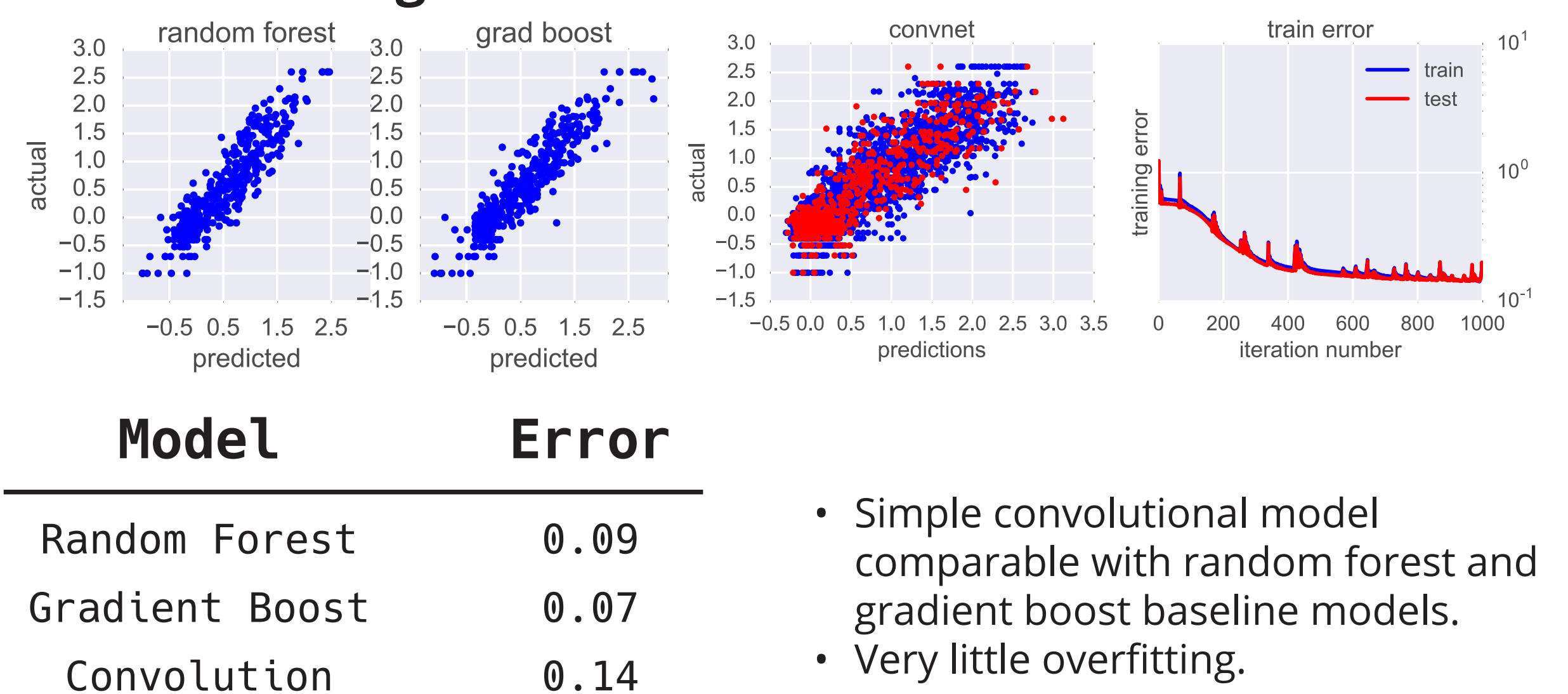
Results

Synthetic Data

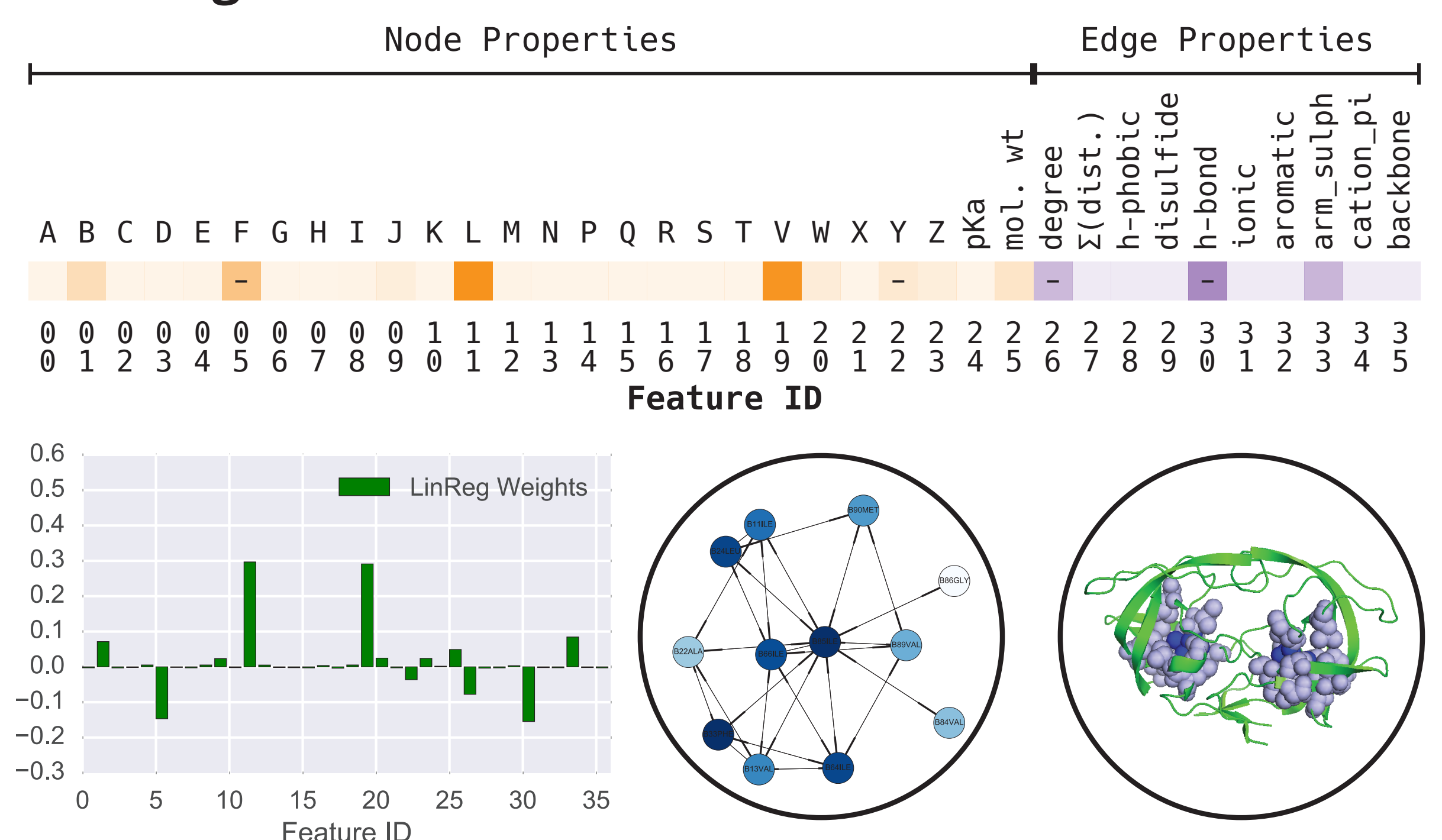


- Learn mathematical transforms on top of integer graphs.
- Deeper networks converge in fewer iterations with smaller error.
- Graph-based convolutional neural nets work on simple mathematical functions.

HIV Protease Drug Resistance



Visualizing Learned Features



- Able to recapitulate known mutations involved in FPV drug resistance.
- Interpretable: hydrophobic network of amino acids implicated in FPV resistance.
- (left) Dark nodes: highly activating; light nodes: weakly activating
- (right) Green ribbon: backbone; Dark blue spheres: top activating nodes; Light blue spheres: neighbors.

Future Work

- Learning Capacity: neural network architecture improvements; prevent overfitting.
- Interpretability: better visualizations of convolutional feature maps.
- Applications: pathogen genomic surveillance, chemical surveillance.