# Chemical Property Prediction via Graph Knowledge Transfer

**Zhenbang Wu**[12*], Haonan Wang[2*], Ziniu Hu[3], Yizhou Sun[3]

[1] Zhejiang University
[2] University of Illinois at Urbana-Champaign
[3] University of California, Los Angeles
* Equal Contribution
{zw12, haonan3}@illinois.edu, yzsun@cs.ucla.edu

# CONTENT

- Introduction
- Related Work
  - Graph-Level Classification / Regression
  - Transfer Learning
  - Multi-Task Learning
- Method
  - Hypergraph Knowledge Transfer
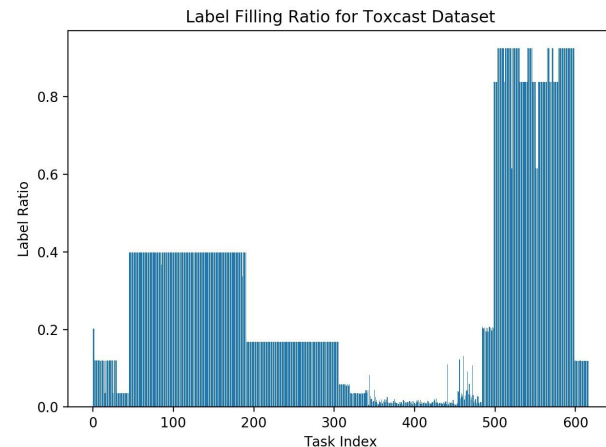- Experiment
- Next Step
- Conclusion

# Introduction

# Introduction

- Predicting molecular properties
  - A fundamental problem in Biomedicine and Chemistry
  - Expensive and time-consuming
- Use of deep learning (DL)
  - Speed-up the process
  - Better predict molecular properties

Mengying Sun, Sendong Zhao, Coryandar Gilvary, Olivier Elemento, Jiayu Zhou, Fei Wang, Graph convolutional networks for computational drug development and discovery, Briefings in Bioinformatics, , bbz042, https://doi.org/10.1093/bib/bbz042

# Introduction

- Practical effect of DL is limited
  - Require large amounts of labeled data
- Usually, in Biomedicine and Chemistry
  - Fully labeling a dataset is unaffordable [1]
  - Label ratio between properties is imbalanced [2]
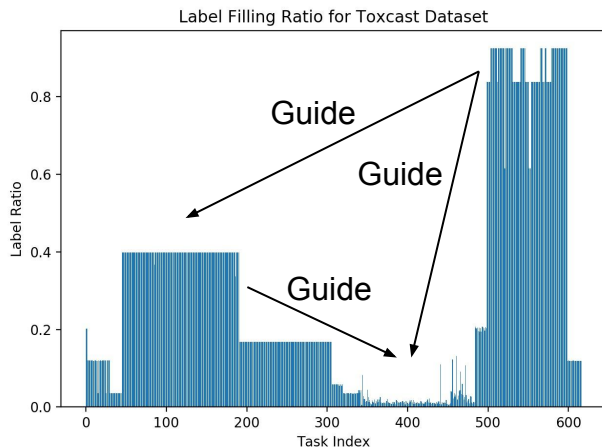


Label Filling Ratio for Toxcast Dataset

[1] H. Altae-Tran, B. Ramsundar, A. S. Pappu, V. Pande. Low Data Drug Discovery with One-Shot Learning. ACS Cent. Sci. 2017, 3 (4), 283−293.
[2] W. Lin, D. Xu, Imbalanced multi-label learning for identifying antimicrobial peptides and their functional types. Bioinformatics, Volume 32, Issue 24, 15 December 2016, Pages 3745.

# Introduction

- Intuition: leverage task dependency
  - Knowledge extracted from fully labeled properties can enhance the prediction of properties with few labels
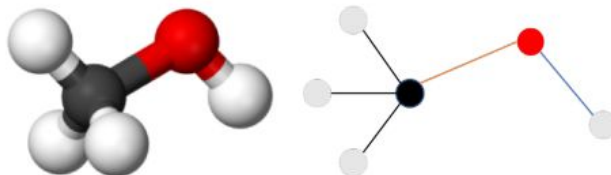


Label Filling Ratio for Toxcast Dataset

| Property | ESRE BLA | APR HepG2 |
|---|---|---|
| Mol 1 | 1 | x |
| Mol 2 | 0 | 0 |
| Mol 3 | 0 | x |
| Mol 4 | 1 | 1 |
| Mol 5 | 1 | x |
| … | … | … |
| Mol 8597 | 0 | x |
| Mol 8598 | 1 | x |
| Label Ratio | 0.84 | 0.12 |

ToxCast dataset. x: no label.

Guide

# Related Work:
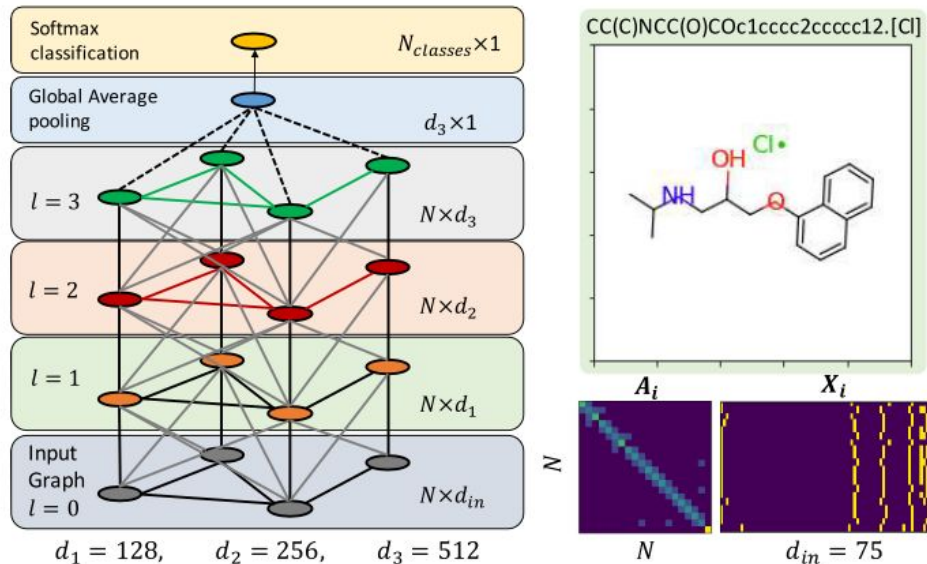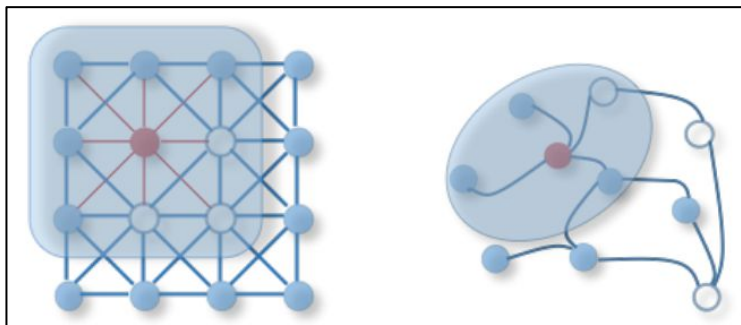# Graph-Level Classification / Regression

# Graph-Level Classification / Regression

- Molecule graph
  - Atoms -> nodes
  - Chemical bonds -> edges



- Molecule -> Graph-Level Prediction
  - e.g. toxicity, solubility, side effect

# Graph Classification / Regression Model



## Suffer from the lack of labeled data

Z. Wu *et al.*, "A Comprehensive Survey on Graph Neural Networks."
Pope, P.E.; Kolouri, S.; Rostami, M.; Martin, C.E.; Hoffmann, H. Explainability Methods for Graph Convolutional Neural Networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, 16–20 June 2019.
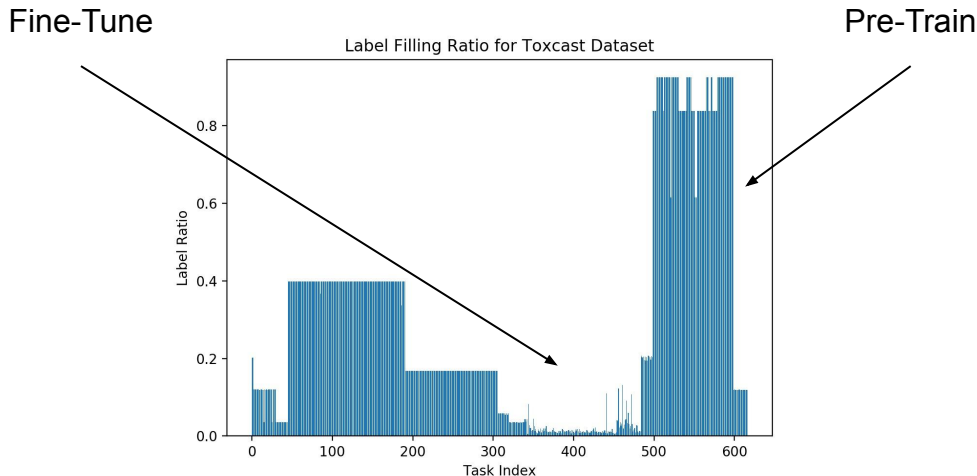
# Related Work: Transfer Learning

# Transfer Learning

- Pre-train the model on properties with abundant labels
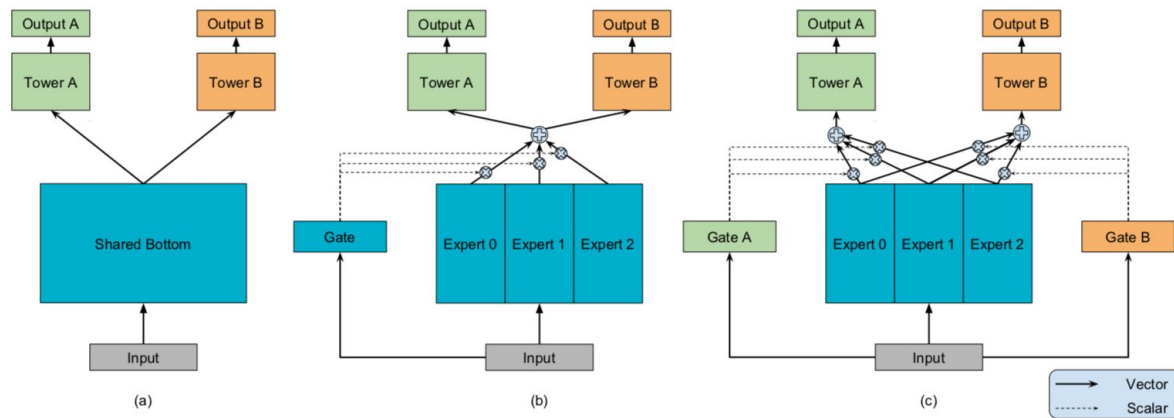- Fine-tune the model on properties with few labels



Fine-Tune

Pre-Train

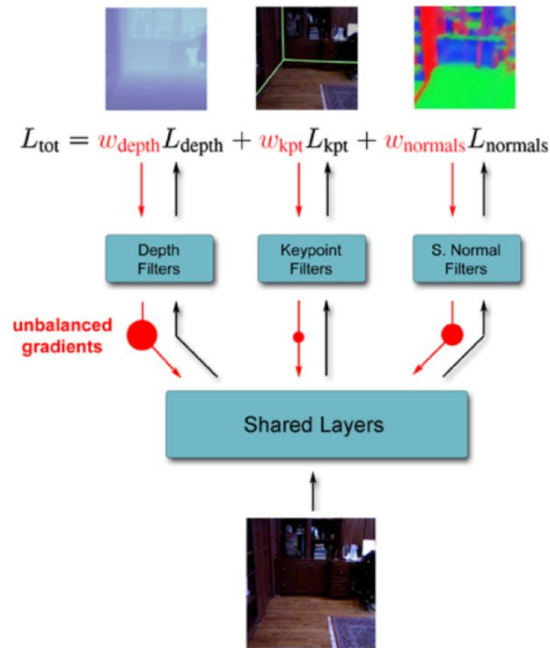Not End-to-End
Not Time-Efficient

# Related Work: Multi-Task Learning

# Multi-Task Learning

- Train multiple tasks together
  - Time efficient
  - Leverate knowledge among tasks



Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. 2018. Modeling Task Relationships in Multi-task Learning with Multi-gate Mixture-of-Experts. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '18). ACM, New York, NY, USA, 1930-1939. DOI: https://doi.org/10.1145/3219819.3220007
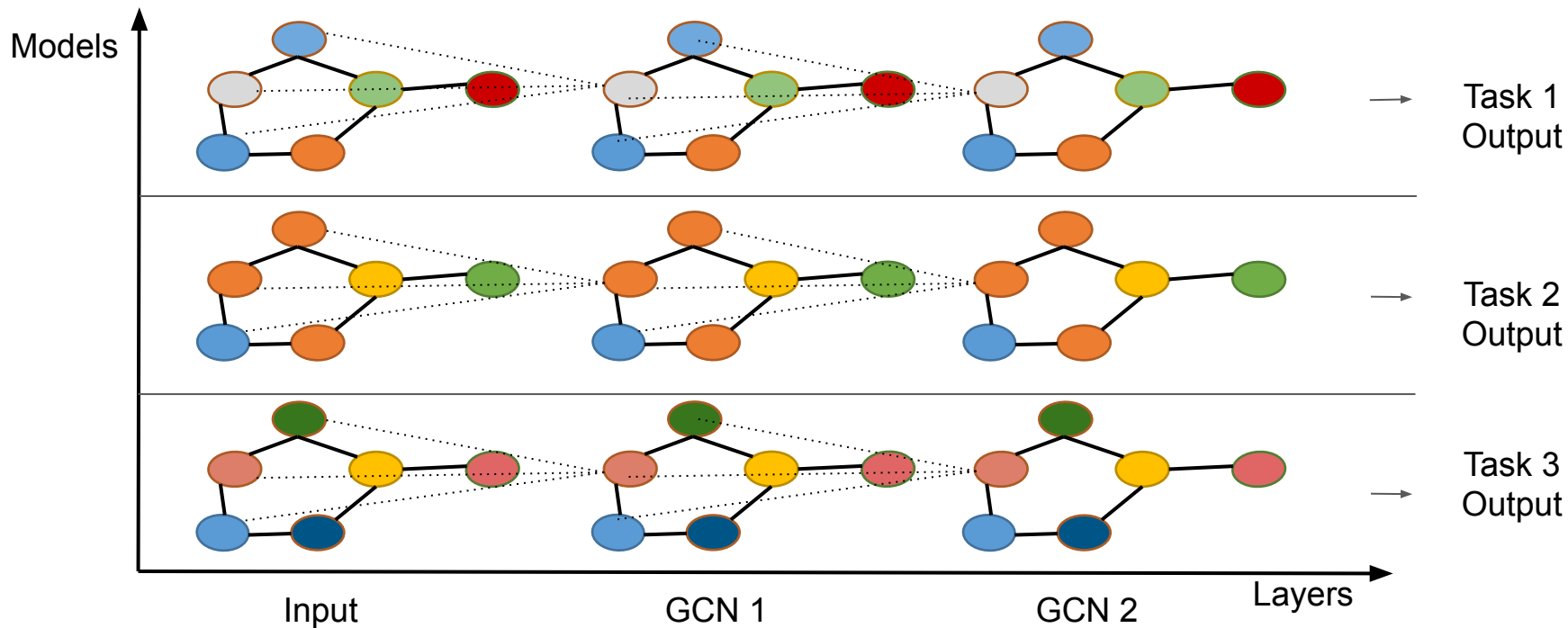
# Multi-Task Learning

- Imbalanced label ratio
- Imbalanced gradients
  - Tasks with abundant labels / larger gradients will dominate the model
- Neglect interaction among different tasks



$$L_{tot} = w_{depth}L_{depth} + w_{kpt}L_{kpt} + w_{normals}L_{normals}$$

Zhao Chen Vijay Badrinarayanan Chen-Yu Lee Andrew Rabinovich GradNorm: Gradient Normalization for Adaptive Loss Balancing in Deep Multitask Networks. 2017 abs/1711.02257 CoRR
http://arxiv.org/abs/1711.02257 db/journals/corr/corr1711.html#abs-1711-02257

# Method: Hypergraph Knowledge Transfer

# Normal GCN

# Hypergraph GCN



Models

Task 1 Output

Task 2 Output

Task 3 Output

Input    GCN 1    GCN 2    Layers
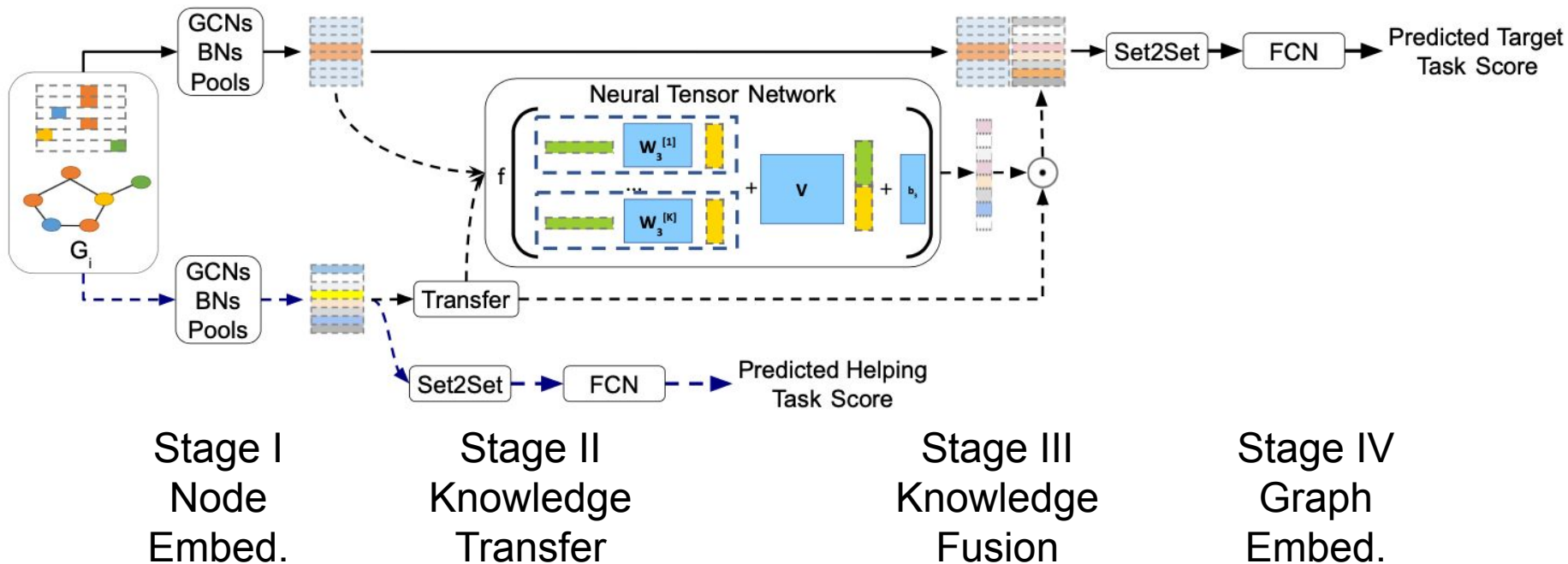
Data-Dependent Task Dependency Graph

# Hypergraph Knowledge Transfer

- Each task will have its own base model
- Calculate the data-dependent task dependency graph
- Aggregate representation from different task specific models

- Mutually enhance performance on all tasks
- Mining the task-level dependency

# Base Model (Dual-Task)



Stage I
Node
Embed.

Stage II
Knowledge
Transfer

Stage III
Knowledge
Fusion

Stage IV
Graph
Embed.

# Experiment

# Experiment Setting

| Dataset | Graph Meaning | #Graphs | #Tasks |
|---------|---------------|---------|--------|
| TOX21 | Qualitative Toxicity Measurements | 7831 | 12 |
| SIDER | Adverse Drug Reactions | 1427 | 27 |

Target tasks (10% training):
- SR-ARE (TOX21)
- Investigations (SIDER)

Helping tasks: (90% training)
- SR-MMP (TOX21)
- Vascular Disorders (SIDER)

# Experiment Results

| Model | TOX21 | | SIDER | |
| --- | --- | --- | --- | --- |
| | Target Task (10% Training) | Helping Task (90% Training) | Target Task (10% Training) | Helping Task (90% Training) |
| Single-task Model | | | | |
| GCN [3] | 0.6776 | 0.8638 | 0.5938 | 0.6266 |
| MoleculeNet [4] | 0.7156 | 0.8315 | 0.6189 | 0.6294 |
| Our | 0.7385 | 0.9096 | 0.6266 | **0.8212** |
| Multi-task Model | | | | |
| MoleculeNet [4] | 0.7298 | 0.8382 | 0.6315 | 0.6503 |
| Our | **0.7762** | **0.9233** | **0.6569** | 0.8037 |

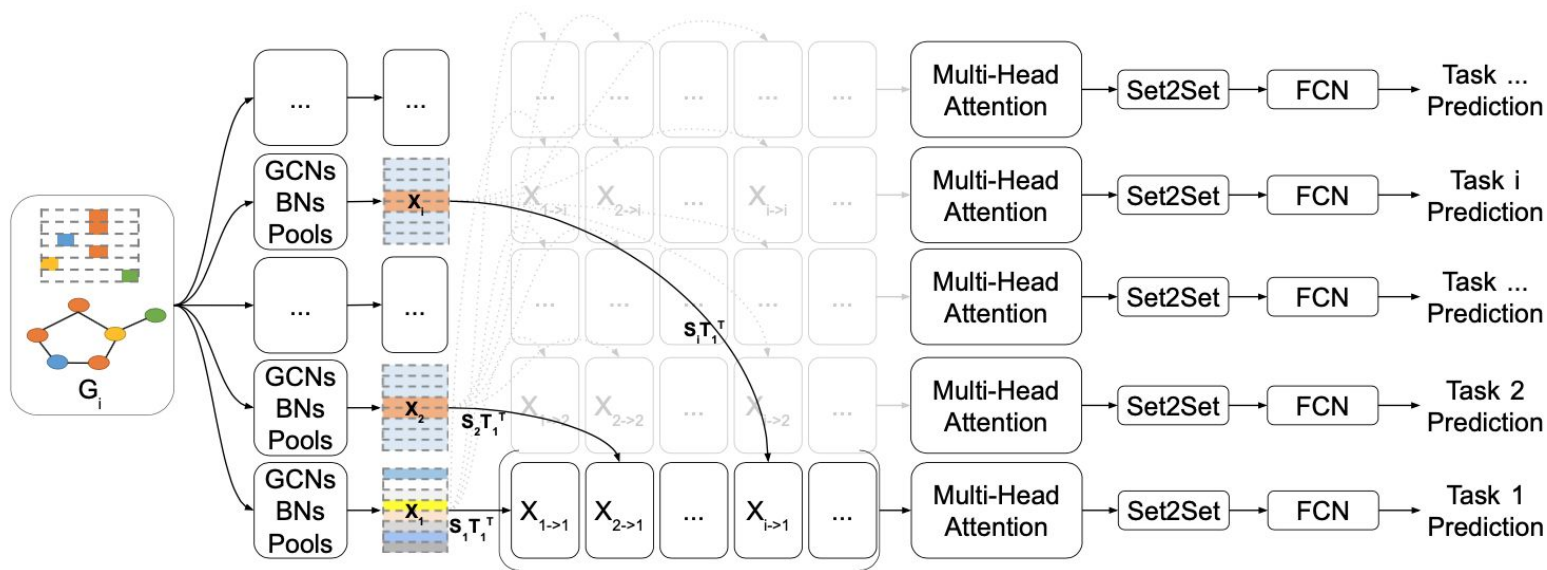- Score will decrease if we lower the ratio of training data
- Score will increase if we leverage knowledge between tasks

Z. Wu, B. Ramsundar, E. N. Feinberg; J. Gomes,C. Geniesse, A. S. Pappu, K. Leswing, V. Pande. MoleculeNet: A Benchmark for Molecular Machine Learning, 2017.
https://arxiv.org/abs/1703.00564

# Next Step

# Backbone Model (Multi-Task)

- Problems of dual-task model
  - For every task pair (i, j), We need a transfer module $f_{i \rightarrow j}$
    - $O(k^2)$, k is #tasks -> not scalable
  - Ignore the tasks relation at a higher level
- Insight:
  - Decompose the transfer module $f_{i \rightarrow j} = S_i T_j^{'}$
    - Each task i only need to store $S_i$ and $T_j$
    - $O(k)$, k is #tasks
  - Explicitly model task-level and graph-level relations
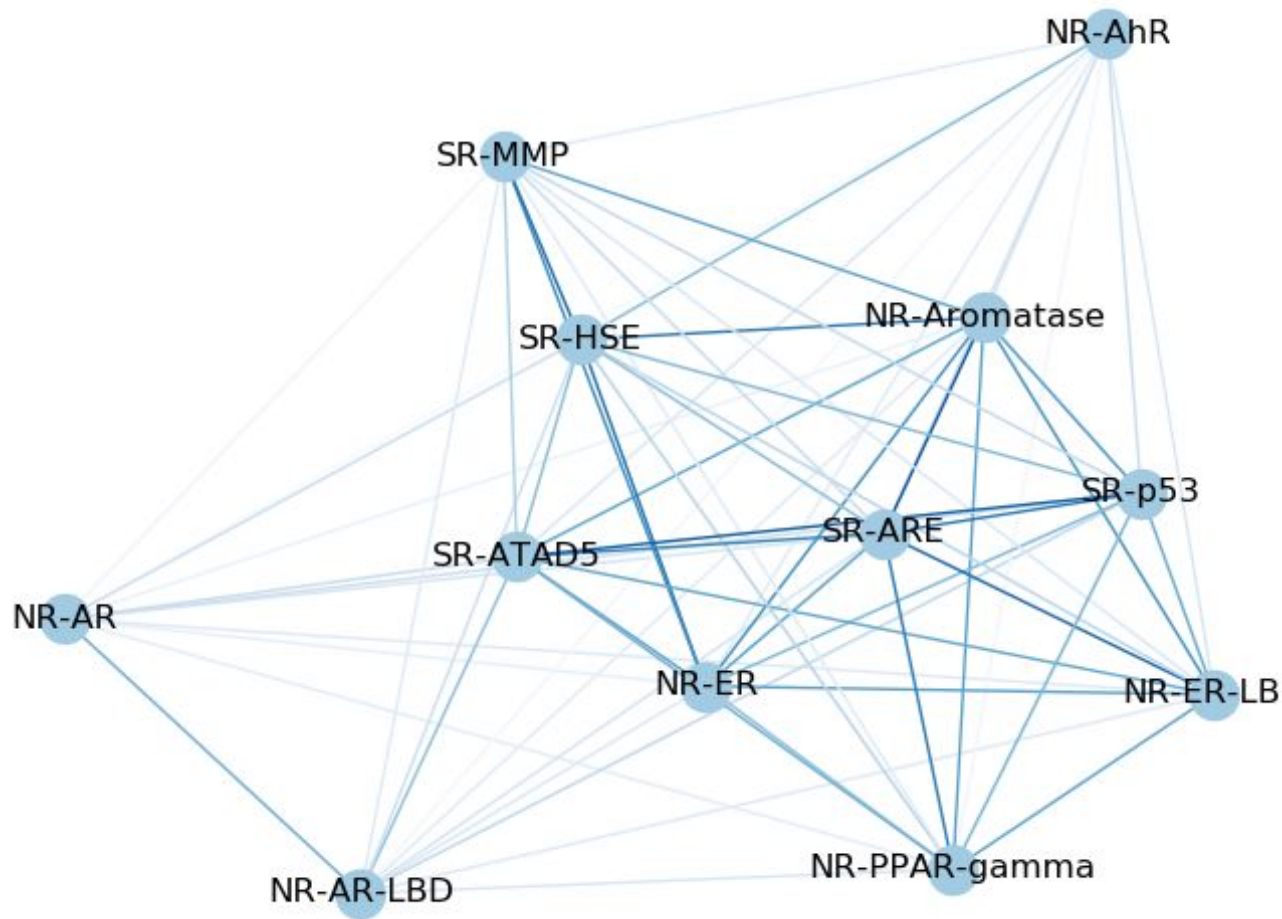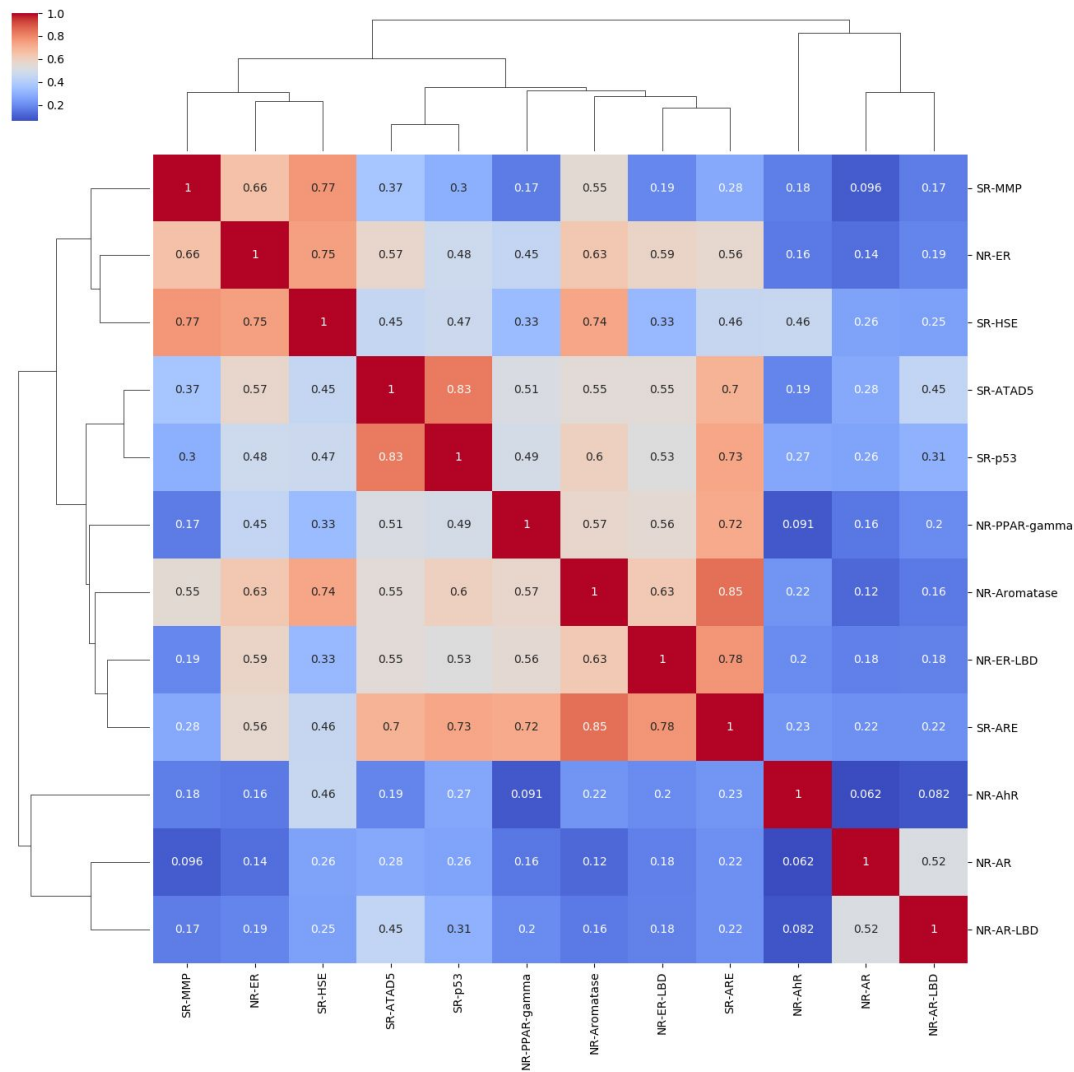
# Backbone Model (Multi-Task)



$$\hat{\boldsymbol{X}}_j = \Sigma_{i=1}^k [\mathrm{Softmax}(D_{i \to j} \boldsymbol{A}_{i \to j}) \cdot \boldsymbol{X}_{i \to j}]$$

$$\boldsymbol{X}_{i \to j} = \boldsymbol{X}_i \boldsymbol{S}_i \boldsymbol{T}_j^T \text{ , where } \boldsymbol{S}_i, \boldsymbol{T}_j \in R^{d \times d'}, \boldsymbol{X}_{i \to j} \in R^{n \times d}$$

$$\boldsymbol{A}_{i \to j} = \boldsymbol{X}_{i \to j} \boldsymbol{Q}_i (\boldsymbol{X}_j \boldsymbol{K}_j)^T, \text{ where } \boldsymbol{Q}_i, \boldsymbol{K}_j \in R^{d \times d''}, \boldsymbol{A}_{i \to j} \in R^{n \times n}$$

# Some Results

Dependency Between Target and Source Tasks

# Conclusion

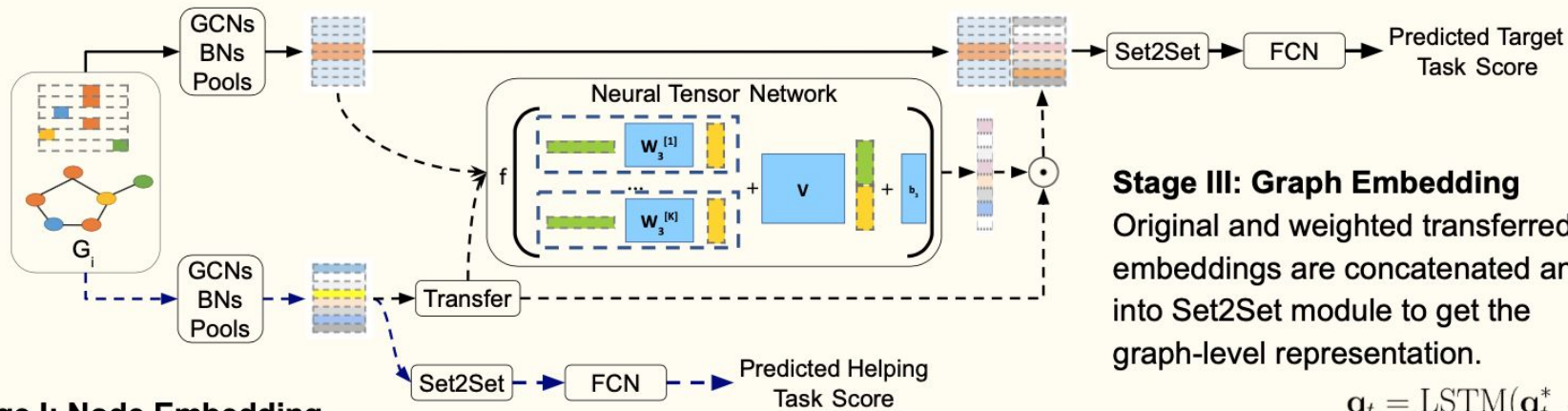# Conclusion

- Lack of labeled data and imbalanced label ratio limits the effect of DL in Chemistry and Biomedicine
- Contribution of our work: transfer + fusion
  - Introduce hypergraph to transfer knowledge between properties
  - Introduce novel attention mechanism to fuse transferred knowledge
  - Explore the hidden dependency structure between tasks
  - Improve dual-task's AUC-ROC score by 6.9%

zw12@illinois.edu

# Q&A

- Introduction
- Related Work
  - Graph-Level Classification / Regression
  - Transfer Learning
  - Multi-Task Learning
- Method
  - Hypergraph Knowledge Transfer
- Experiment
- Next Step
- Conclusion

# MODEL DETAIL



**Stage I: Node Embedding**

$$\text{Conv}(A, X) = \hat{D}^{-1/2} \hat{A} \hat{D}^{-1/2} X \Theta$$

$$\text{BN}(X) = \frac{X - \mathbb{E}[X]}{\sqrt{\text{Var}[X] + \epsilon}} * \gamma + \beta$$

$$\text{Pool}(v) = \max\{\max_{(u,v) \in E}\{u, v\}\}$$

**Stage II: Knowledge Transfer**

Two linear layers with ReLU activation transfer the node embeddings from helping task to target task.

Neural Tensor Network models the node-level interaction and decides the transfer weights.

**Stage III: Graph Embedding**

Original and weighted transferred embeddings are concatenated and fed into Set2Set module to get the graph-level representation.

$$\mathbf{q}_t = \text{LSTM}(\mathbf{q}_{t-1}^*)$$

$$\alpha_{i,t} = \text{softmax}(\mathbf{x}_i \cdot \mathbf{q}_t)$$

$$\mathbf{r}_t = \sum_{i=1}^{N} \alpha_{i,t} \mathbf{x}_i$$

$$\mathbf{q}_t^* = \mathbf{q}_t \| \mathbf{r}_t,$$

33

Assume $T_j$ is the target task.

**Transfer the embeddings from $T_i$ to $T_j$ ($i \neq j$):**

$X_{i \to j} = X_i W_{Si} W_{Tj}^T$, where $W_{Si}, W_{Tj} \in R^{d \times d'}$, $X_{i \to j} \in R^{n \times d}$.

**Calculate the node-level attention:**

$\mathrm{ATT}_{i \to j} = X_{i \to j} W_{Qi} (X_j W_{Kj})^T$, where $W_{Qi}, W_{Kj} \in R^{d \times d''}$, $\mathrm{ATT}_{i \to j} \in R^{n \times n}$.

$\widetilde{\mathrm{ATT}}_{i \to j} = \mathrm{Softmax}(\mathrm{ATT}_{i \to j}, \dim = -1) \in R^{n \times n}$.

**Combine embeddings from all tasks w.r.t. node-level attention**

$\hat{X}_{i \to j} = \widetilde{\mathrm{ATT}}_{i \to j} \cdot X_{i \to j}$ where $\hat{X}_{i \to j} \in R^{n \times d}$.

$X_j^{\mathrm{comb}} = \mathrm{cat}([X_j, \hat{X}_{: \to j}]) \in R^{k \times n \times d}$.

$\widetilde{X_j^{\mathrm{comb}}} = \mathrm{Norm}(X_j^{\mathrm{comb}}) \in R^{k \times n \times d}$.

**Merge the embeddings w.r.t. task dependency:**

$\hat{X}_j = W_{Dj} \widetilde{X_j^{\mathrm{comb}}}$, where $W_{Dj} \in R^{1 \times k}$, $\hat{X}_j \in R^{n \times d}$.

**Finally, calculate graph-level embedding:**

$G_j = \mathrm{Readout}(\hat{X}_j) \in R^{1 \times d}$.

# Tox21 Data Challenge

**Training Data** ~12,000 Compounds

**Evaluation Data** ~650 Compounds

**Nuclear Receptor Panel (biomolecular targets)**
- ER-LBD: estrogen receptor alpha, luciferase
- ER: estrogen receptor alpha
- aromatase
- AhR: aryl hydrocarbon receptor
- AR: androgen receptor
- AR-LBD: androgen receptor, luciferase
- PPAR: peroxisome proliferator-activated receptor gamma

**Stress Response Panel**
- ARE: nuclear factor (erythroid-derived 2)-like 2 antioxidant responsive element
- HSE: heat shock factor response element
- ATAD5: genotoxicity indicated by ATAD5
- MMP: mitochondrial membrane potential
- p53: DNA damage p53 pathway