

# Property-Aware Relation Networks for Few-shot Molecular Property Prediction

Dr. Quanming Yao

*Assistant professor, EE Tsinghua*

[qyaoaa@tsinghua.edu.cn](mailto:qyaoaa@tsinghua.edu.cn)

Joint work with [Yaqing Wang, Dejing Dou] (Baidu), [Abulikemu Abuduweili] (CMU)

NeurIPS 2021 (Spotlight)

Code: <https://github.com/tata1661/PAR-NeurIPS21>

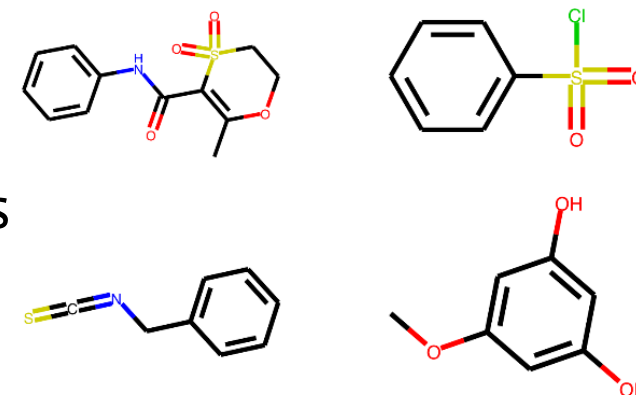


# Outline

- Background
  - Molecular property prediction (MPP)
  - Few-shot learning (FSL)
- Existing works
- The proposed approach
- Summary

# Molecular Property Prediction (MPP)

- Molecules:
  - Mainly micromolecule organics
  - Graph-structured data, Graph Neural Network (GNN) is useful in obtaining its representations
- Properties:
  - Physiology or Toxicity
  - Examples in SIDER :
    - ‘SIDER’ : [ ‘Hepatobiliary disorders’ , ‘Infections and infestations’ , ‘Neoplasms benign, malignant and unspecified (incl cysts and polyps)’ , ... ]



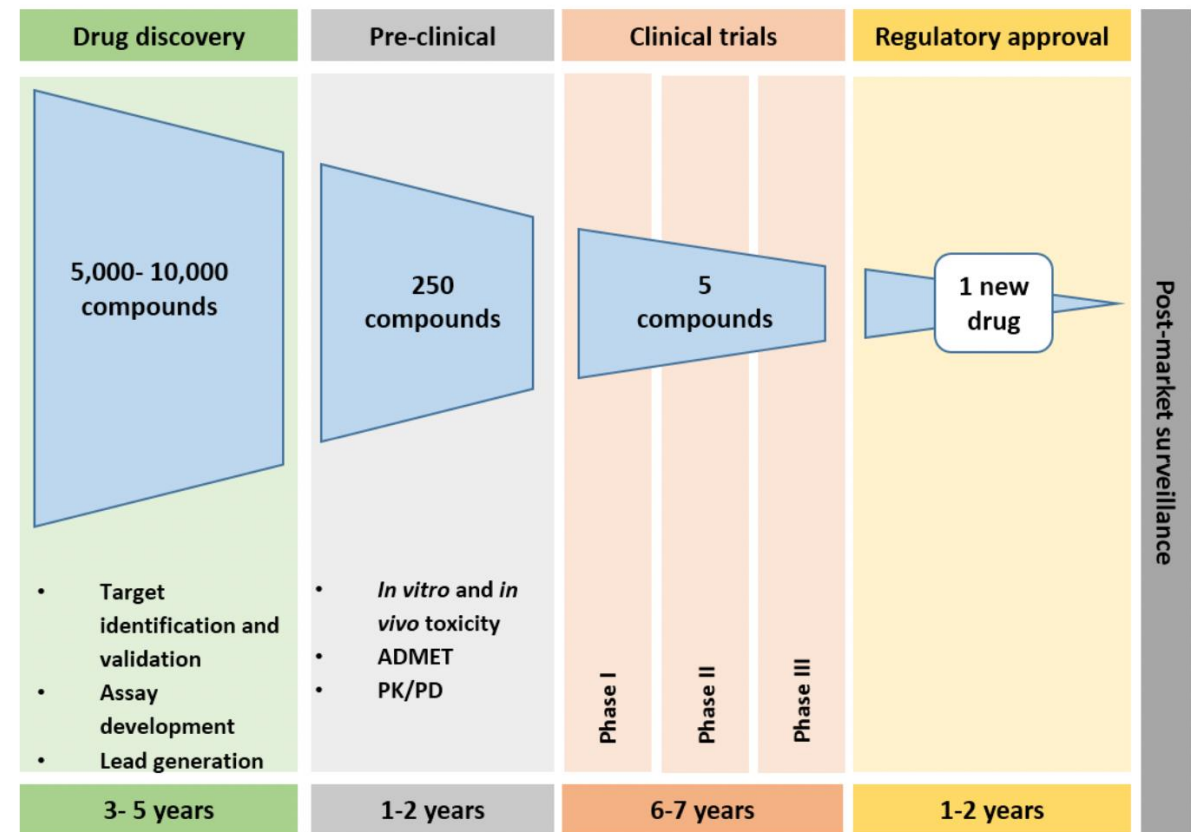
Examples of molecules

# Needs for MPP

Drug discovery targets at finding **new potential** medical **compounds** with **desired properties**

Only a small amount of candidate molecules can **pass virtual screening** to be evaluated in the lead optimization stage

We only have **few molecules** with **known pharmacological properties**



Drug discovery and development timeline from [H. Matthews et al., Proteomes 2016]

# Outline

- Background
  - Molecular property prediction (MPP)
  - Few-shot learning (FSL)
- Existing works
- The proposed approach
- Summary

# Few-shot Learning (FSL)

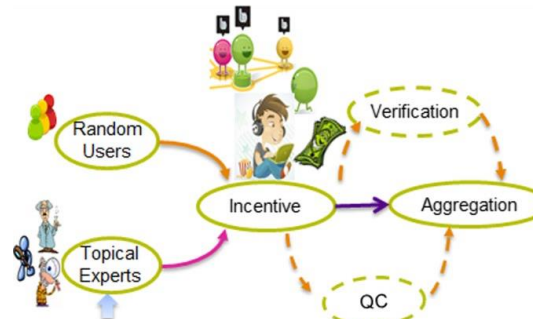
- Definition: A type of machine learning problems contains only a limited number of labeled examples

MPP: intrinsically a few-shot problem

**Shot**: the number of labeled examples

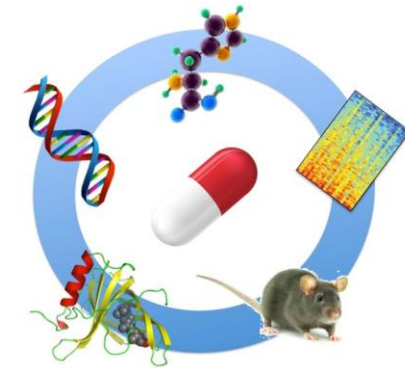
- Typical Scenarios:

Reducing data gathering effort and computational cost



Example: Image / Text Classification  
Labor Intensive / Hard  
(few labeled images/texts)

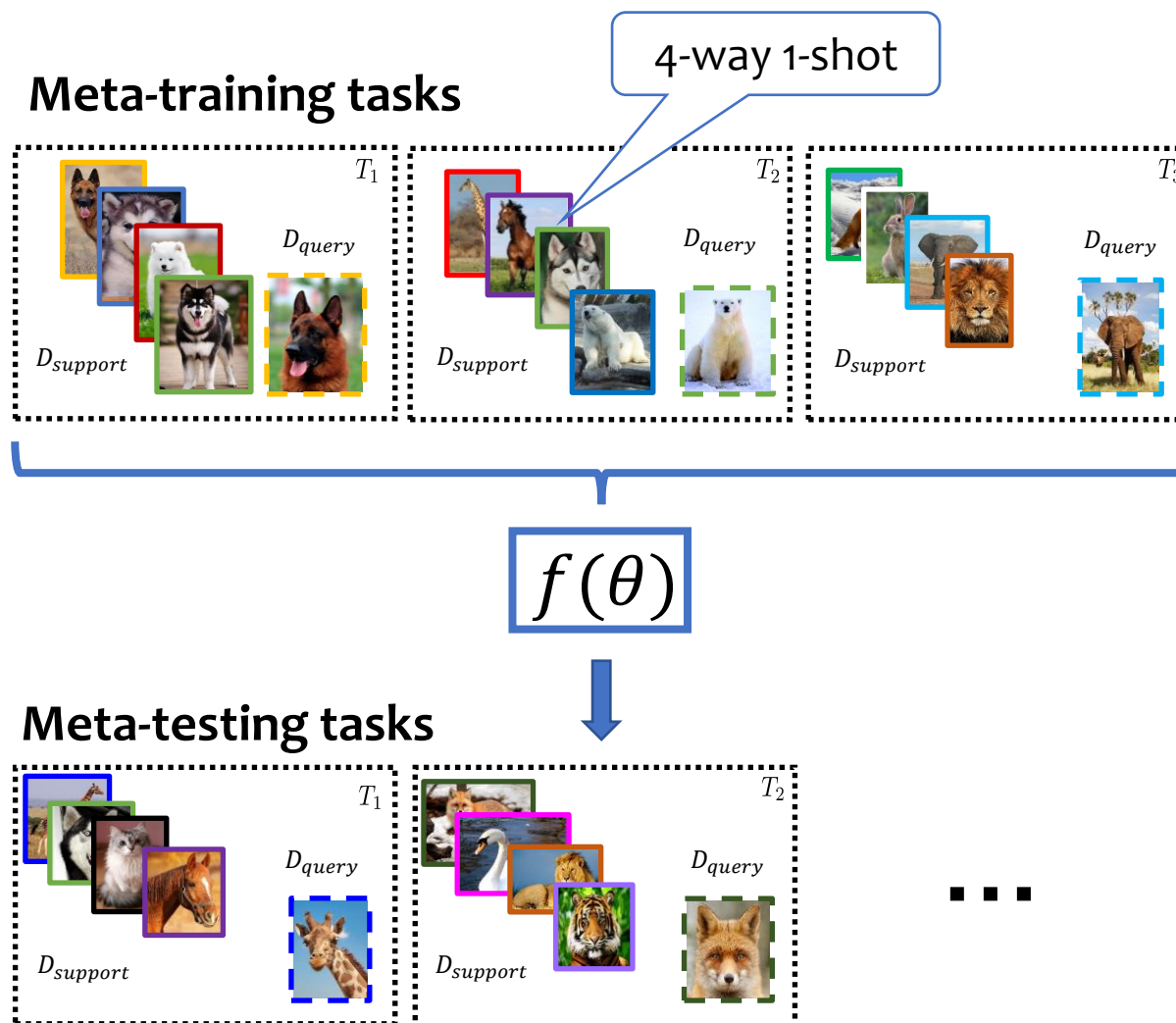
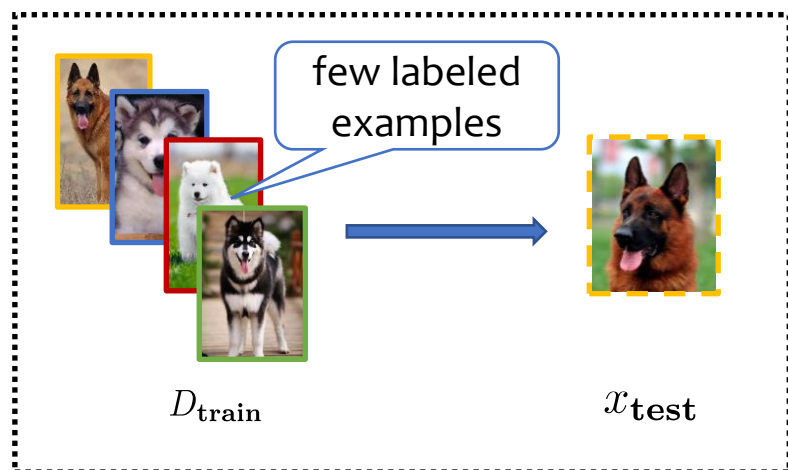
Learning for rare cases



Example: Drug Discovery  
Dangerous / Private / Ethical  
(few labeled drug molecules)

# General Idea

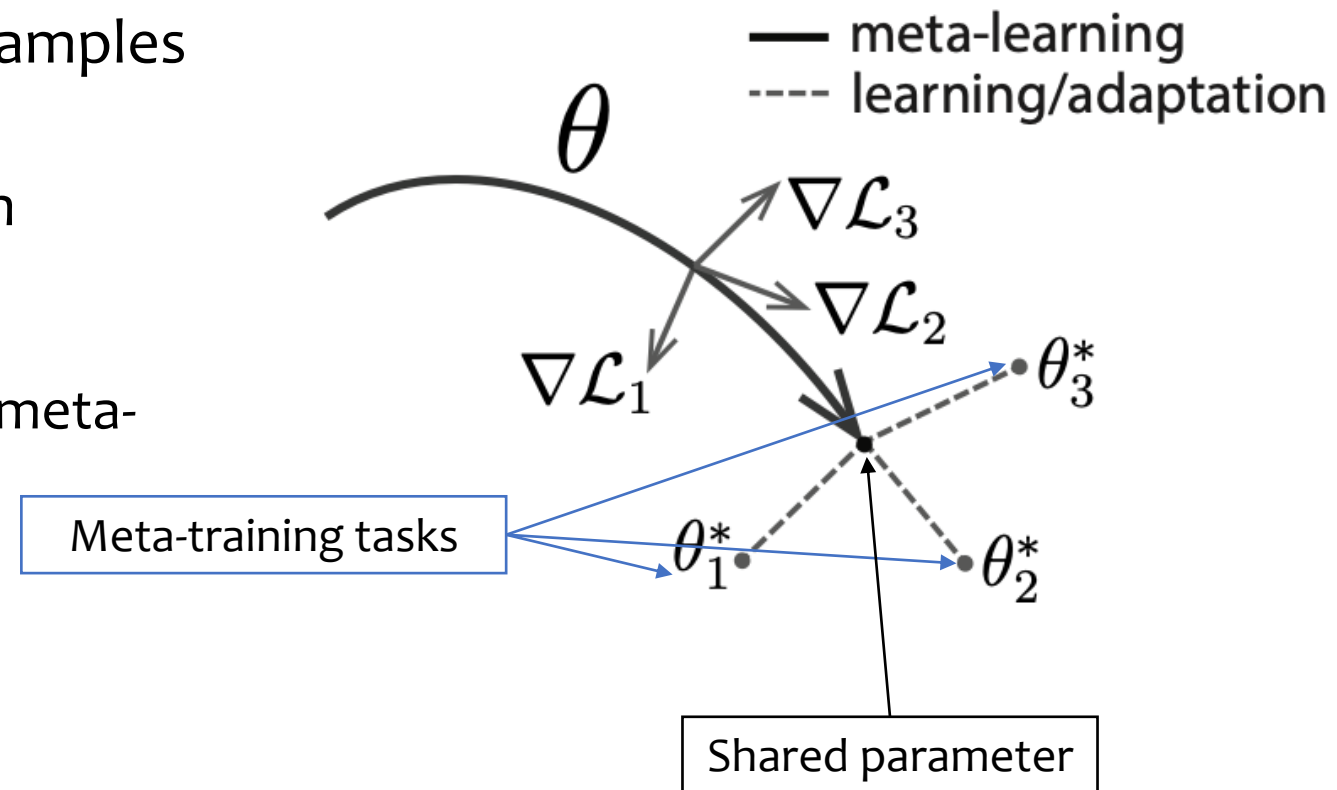
Obtain **prior knowledge** from meta-training (i.e. similar) tasks



# Exemplar Method: MAML

Idea: Train a model on meta-training tasks, such that it can solve new learning tasks using only a small number of training samples

- Prior knowledge: a **shared  $\theta$**  learn from meta-training tasks
- How: using gradient descent to train all meta-training task to get  $\theta$ .





# Exemplar Method: MAML

- Algorithm:

- Step 1: Adapt the model for each meta-training tasks from the **shared  $\theta$**
- Step 2: Update the **share  $\theta$**  by the loss calculated through the models after adaptation

---

**Algorithm 2** MAML for Few-Shot Supervised Learning
 

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

```

1: randomly initialize  $\theta$ 
2: while not done do
3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$ 
4:   for all  $\mathcal{T}_i$  do
5:     Sample  $K$  datapoints  $\mathcal{D} = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$  from  $\mathcal{T}_i$ 
6:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  using  $\mathcal{D}$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation (2)
       or (3)
7:     Compute adapted parameters with gradient descent:
        $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ 
8:     Sample datapoints  $\mathcal{D}'_i = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$  from  $\mathcal{T}_i$  for the
       meta-update
9:   end for
10:  Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  using each  $\mathcal{D}'_i$ 
    and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 2 or 3
11: end while
  
```

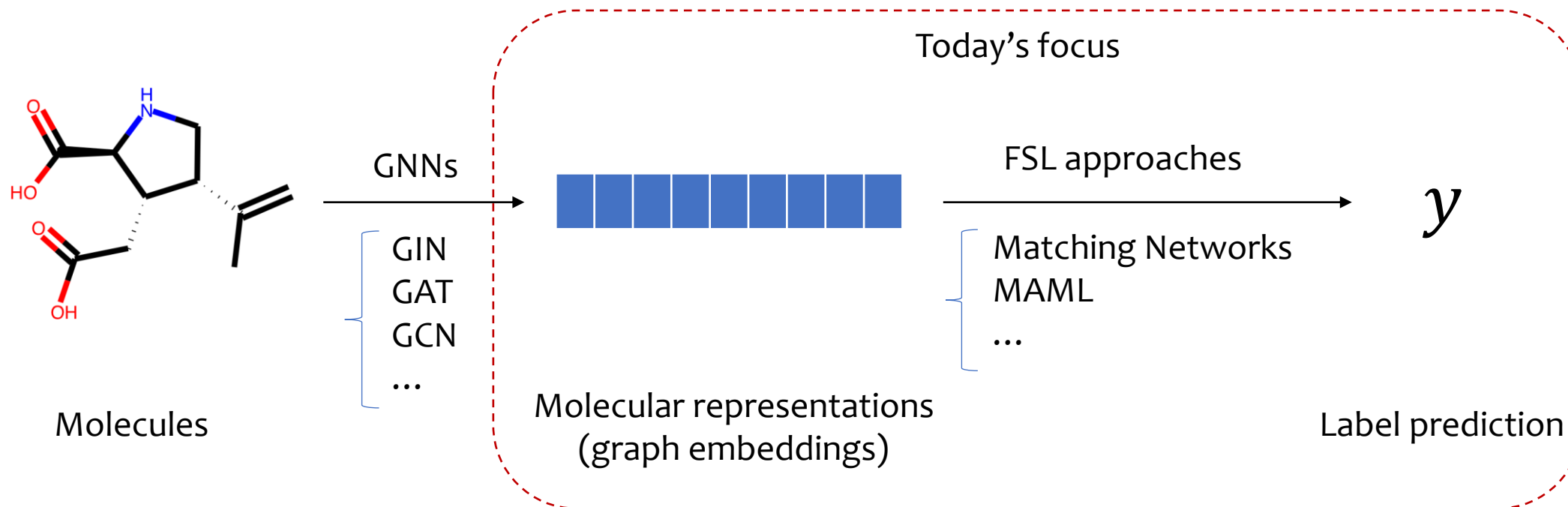
---

Step 1

Step 2

# Take Home Message

- Molecular Property Prediction is a Few-shot graph learning problem



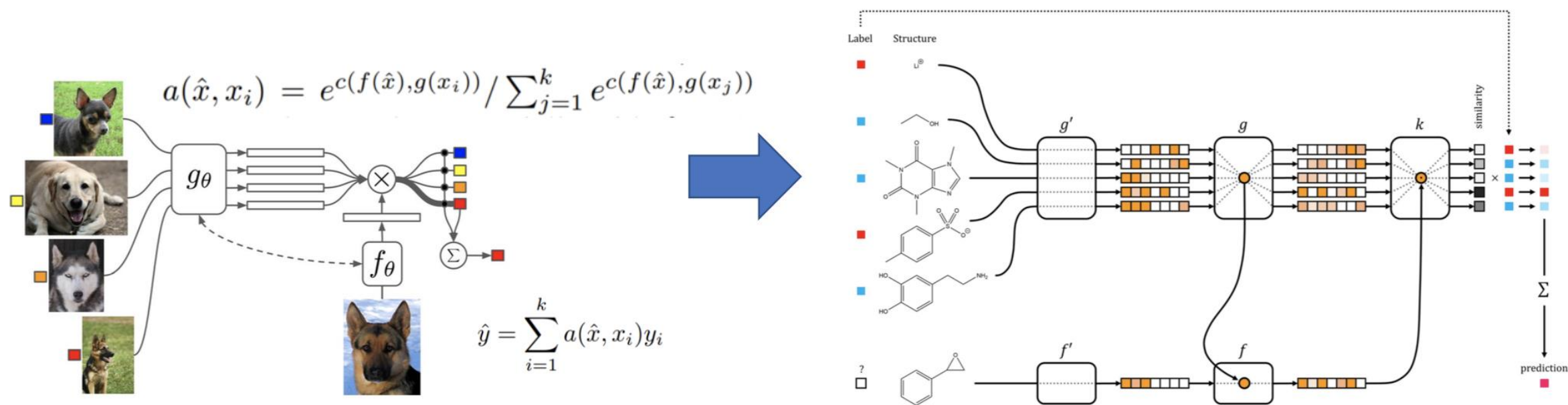
# Outline

- Background
- Existing works
- The proposed approach
- Summary

# Existing Work: IterRefLSTM

## Motivation

- Adapt Matching Networks (one-shot learning) to handle molecular property prediction tasks with few training data
- Propose IterRefLSTM to modify Matching Networks architecture



Oriol et al. 2016. Matching Networks for One Shot Learning, NeurIPS

Han et al. 2017. Low Data Drug Discovery with One-Shot Learning, ACS Central Science

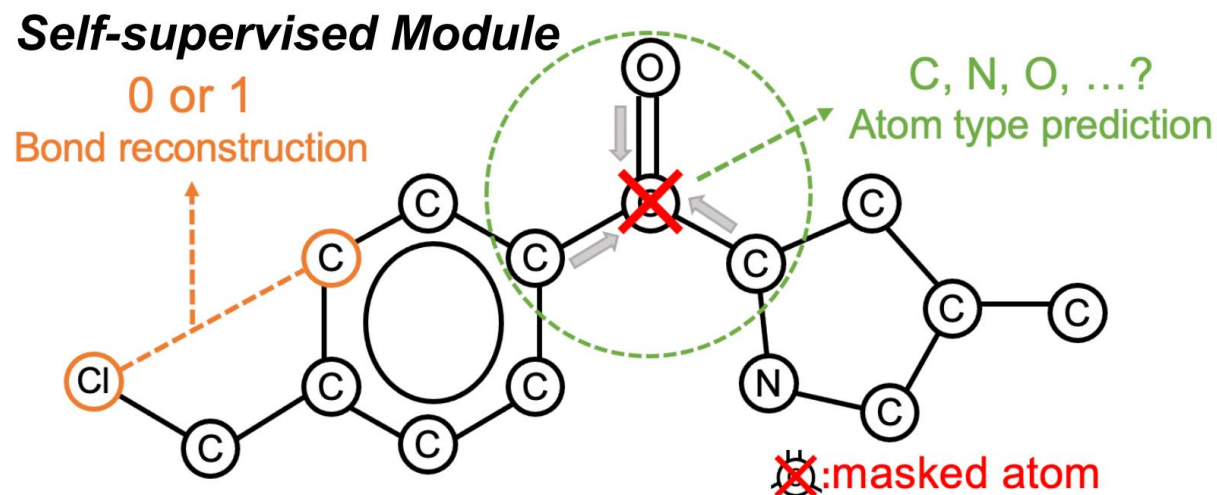
# Existing Work: Meta-MGNN

## Motivation

- Combine MAML and self-supervised learning (i.e., Pre-GNN) for MPP

## Pre-GNN

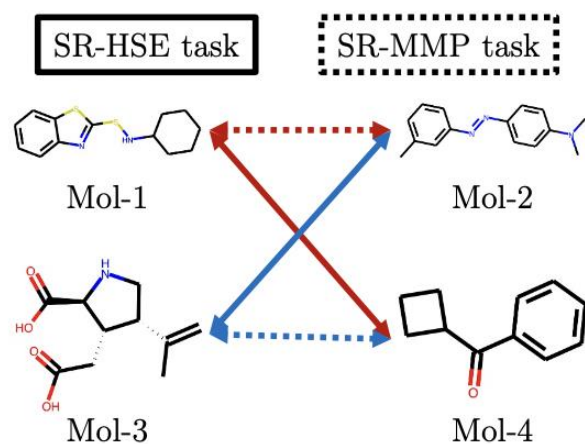
- Exploit the useful unlabeled information in graphs



# Outline

- Background
- Existing works
- The proposed approach
  - Property-Aware Relation networks (PAR)
  - Experiments
- Summary

# Motivation



Molecules		Label	
ID	SMILES	SR-HSE	SR-MMP
Mol-1	<chem>c1ccc2sc(SNC3CCCCC3)nc2c1</chem>	1	1
Mol-2	<chem>Cc1cccc(/N=N/c2ccc(N(C)C)cc2)c1</chem>	0	1
Mol-3	<chem>C=C(C)[C@H]1CN[C@H](C(=O)O)[C@H]1CC(=O)O</chem>	0	0
Mol-4	<chem>O=C(c1ccccc1)C1CCC1</chem>	1	0

Figure 1: Examples of relation graphs for the same molecules coexisting in two tasks of Tox21. Red (blue) edges mean the connected molecules are both active (inactive) on the target property.

Existing works **neglect** two key facts

- The same molecule shows **different properties** in **different MPP tasks**
- **Relationships among molecules** is important and worth-learning

# Property-Aware Relation networks (PAR)

## Summary of ideas

- Property-aware molecular embedding

The same molecule shows **different properties** in **different MPP tasks**

- Relation graph learning

**Relationships among molecules** is important and worth-learning

- Selective update

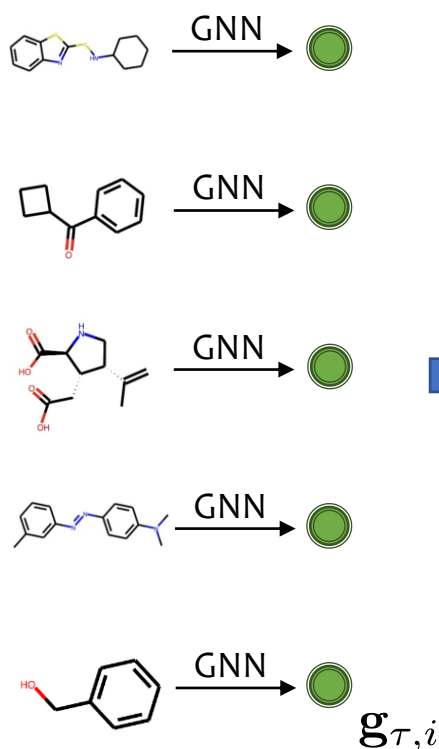
Separately capture the **generic knowledge** shared across different tasks and **property-aware knowledge** within a task



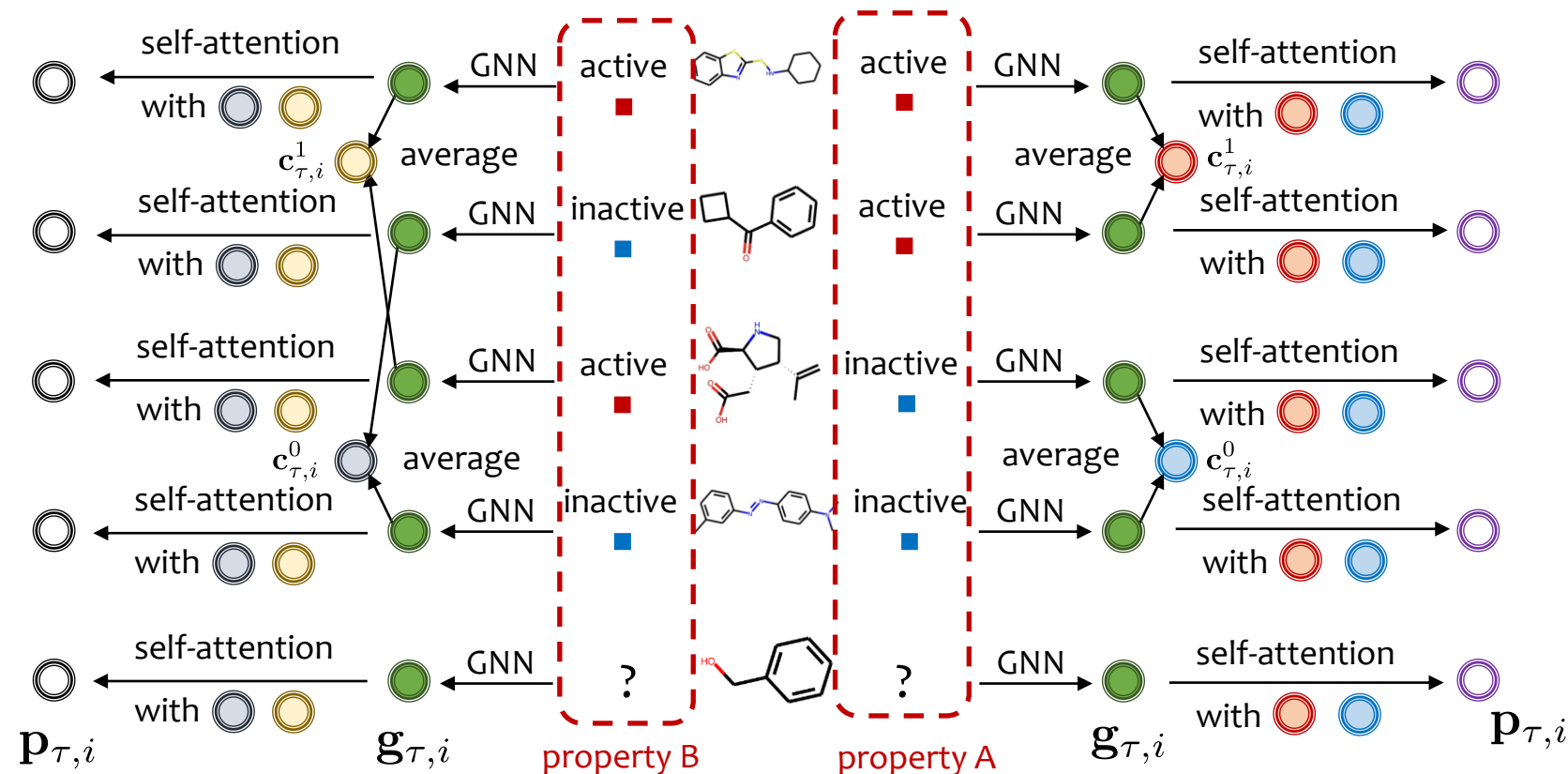
# Property-aware Molecular Embedding

The same molecule shows different properties in different MPP tasks

Classical approaches

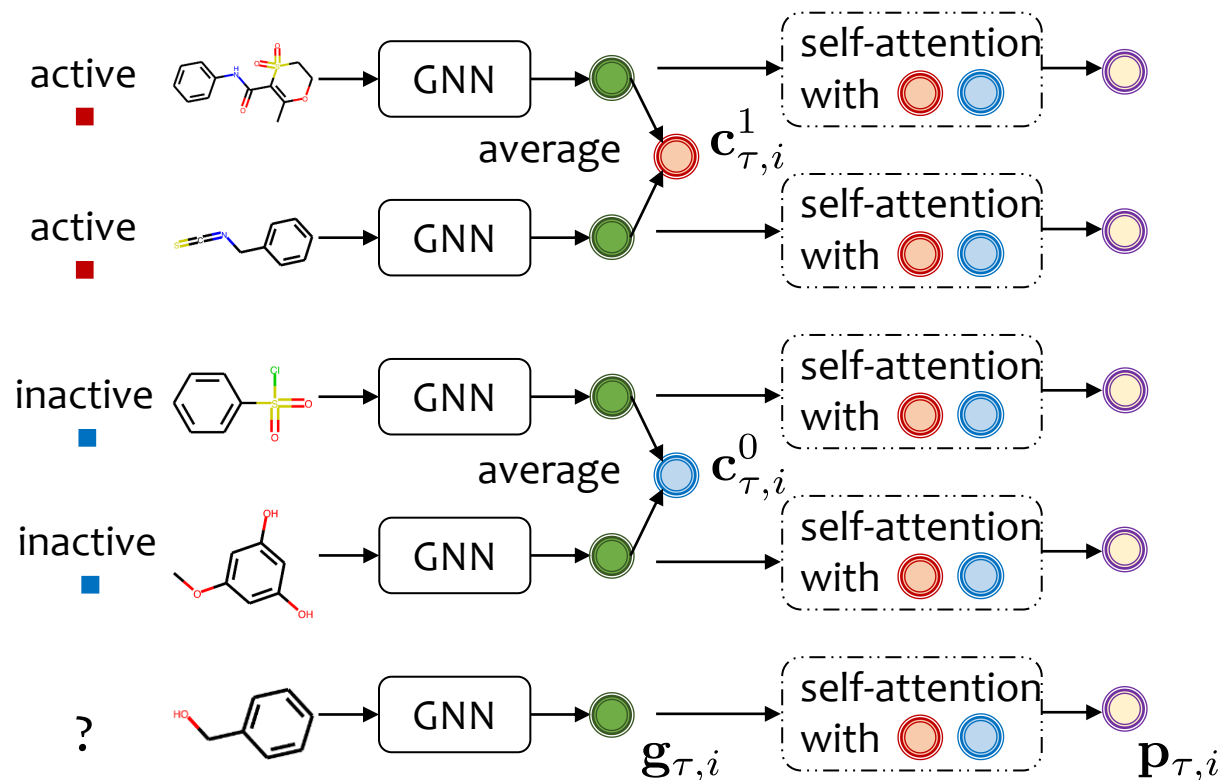


Proposed property-aware Molecular Embedding

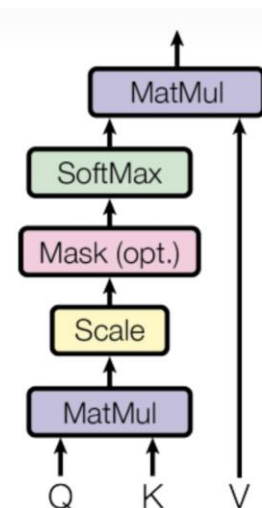


Different labels!

# Property-aware Molecular Embedding



- $g_{\tau,i}$  : molecule representation after GNN
- $c_{\tau,i}^1$  : representative of active molecules
- $c_{\tau,i}^0$  : representative of inactive molecules
- $p_{\tau,i}$  : property-aware embedding
- $Q = K = V = [g_{\tau,i}, c_{\tau,i}^0, c_{\tau,i}^1]$

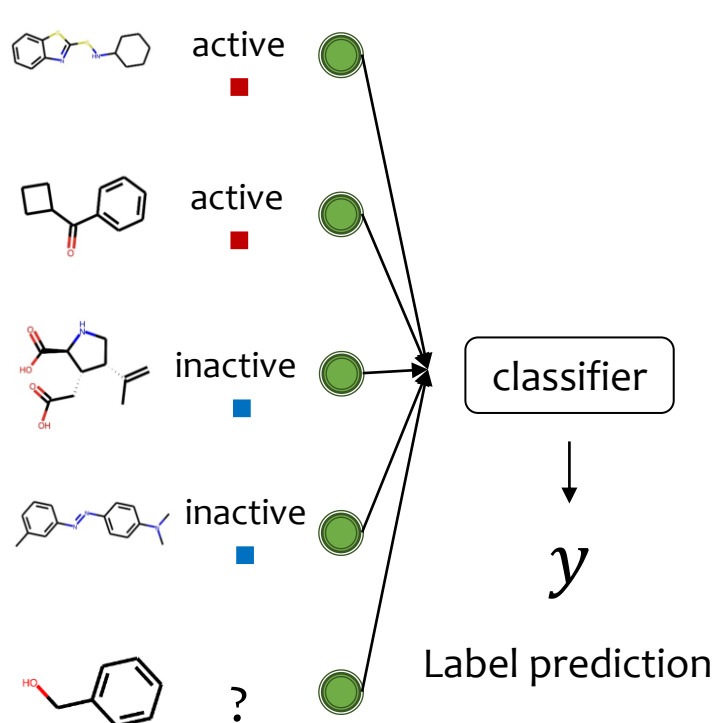


$$\mathbf{b}_{\tau,i} = [\text{softmax}(\mathbf{C}_{\tau,i} \mathbf{C}_{\tau,i}^\top / \sqrt{d_g}) \mathbf{C}_{\tau,i}]_1, \text{ with } \mathbf{C}_{\tau,i}^\top = [g_{\tau,i}, c_{\tau,i}^0, c_{\tau,i}^1] \in \mathbb{R}^{d_g \times 3}$$

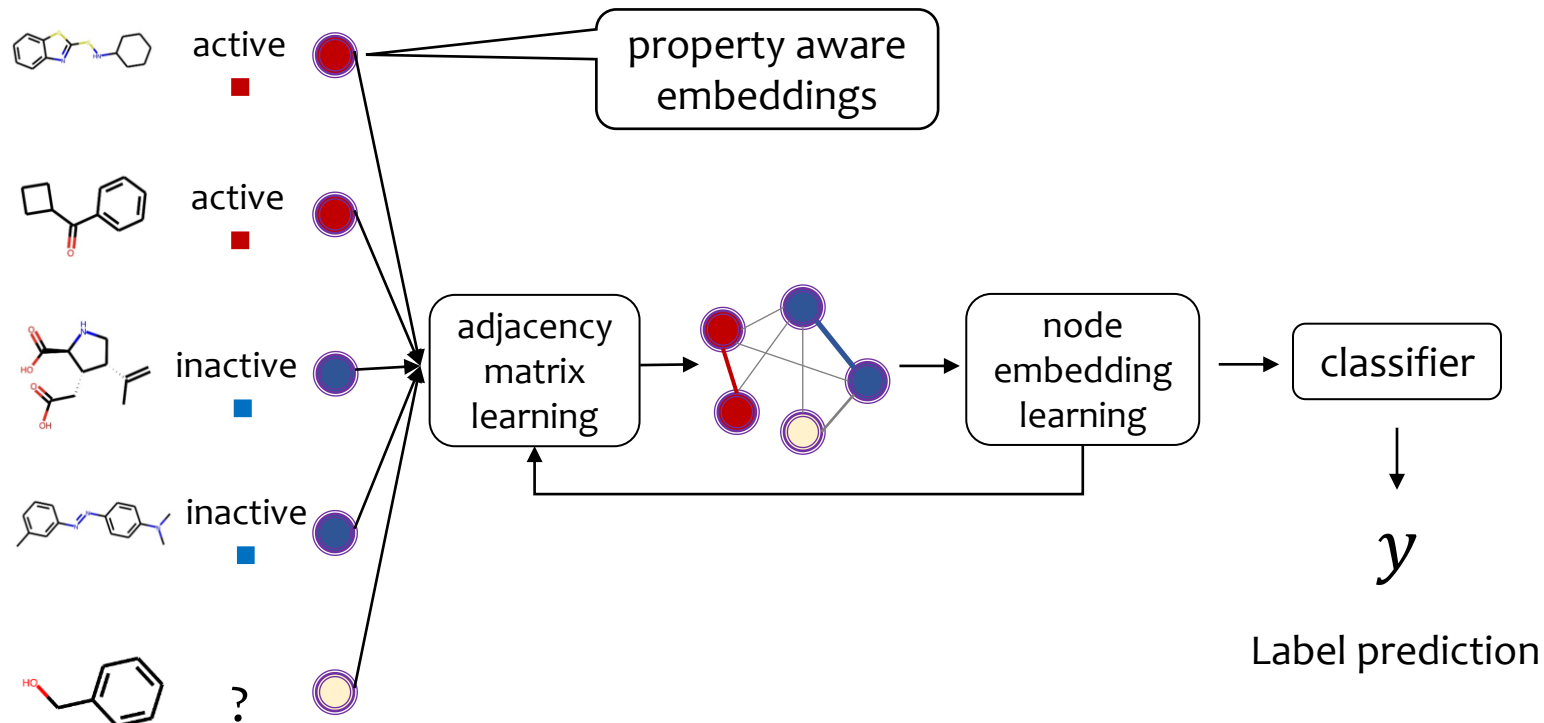
$$\mathbf{p}_{\tau,i} = \text{MLP}_{\mathbf{W}_p}(\text{concat}[g_{\tau,i}, \mathbf{b}_{\tau,i}])$$

# Relation Graph Learning

Relationships among molecules is important and worth-learning



Classical approaches

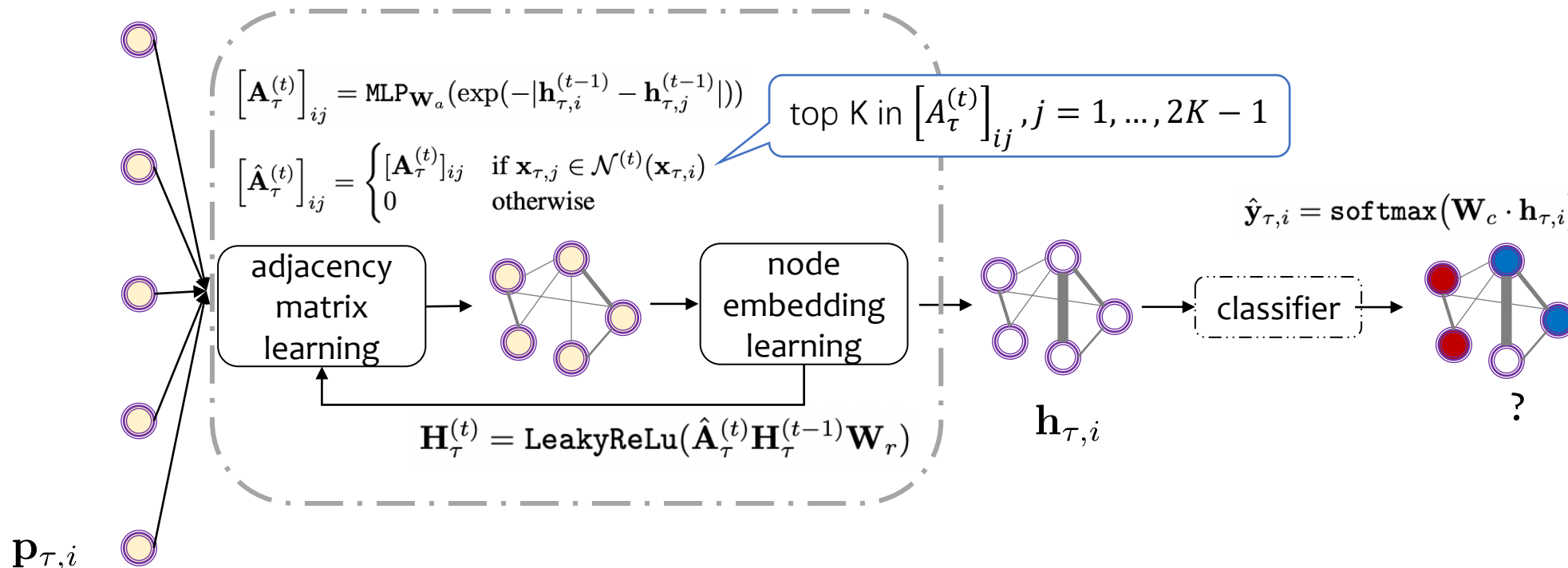


Graph Construction in PAR

# Relation Graph Learning

As Relationships among molecules is important and worth-learning, we

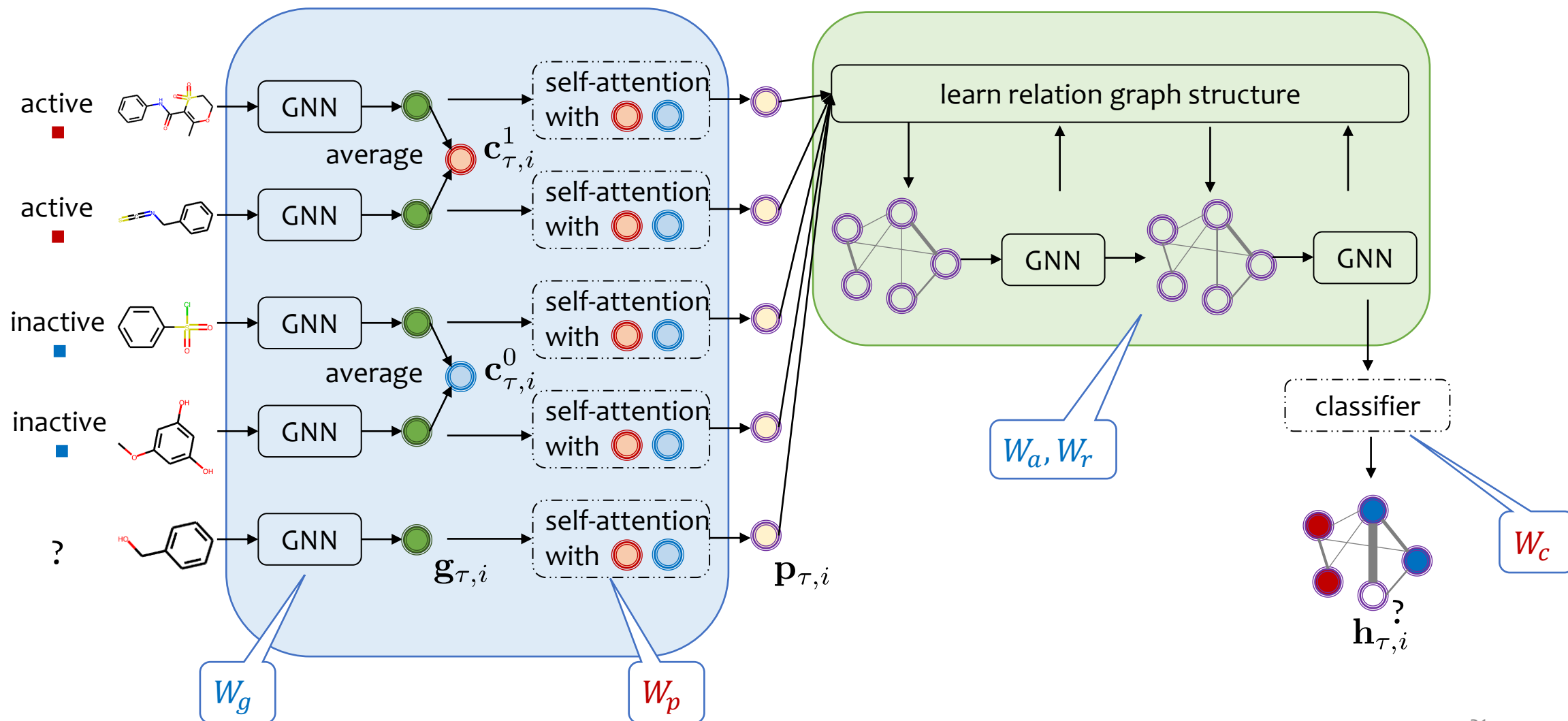
- jointly **estimate molecular relation graph** and **refine molecular embeddings**
- then can propagate limited **labels** efficiently between similar molecules



# PAR Framework

$\theta = \{W_g, W_a, W_r\}$ , generic knowledge, shared across tasks

$\Phi = \{W_p, W_c\}$ , property-aware knowledge, specific to one task



# Training and Inference

Denote PAR as  $f_{\theta, \phi}$

- $\theta = \{W_g, W_a, W_r\}$ : parameters of molecular encoder and relation graph learning module
- $\phi = \{W_p, W_c\}$ : parameters of property-aware embedding function and classifier

We learn from **a set of meta-training tasks** a good shared parameter

$$\min_{\theta, \Phi} \sum_{\tau=1}^{N_t} \mathcal{L}(\mathcal{Q}_{\tau}, f_{\theta, \Phi_{\tau}})$$

Within each task, we **fix**  $\theta$  while **fine-tune**  $\Phi$  as  $\Phi_{\tau}$

Ground-truth labels

$$\mathcal{L}(\mathcal{S}_{\tau}, f_{\theta, \Phi}) = \sum_{(\mathbf{x}_{\tau, i}, y_{\tau, i}) \in \mathcal{S}_{\tau}} -\mathbf{y}_{\tau, i}^{\top} \cdot \log(\hat{\mathbf{y}}_{\tau, i}) + \underbrace{\|[\mathbf{A}_{\tau}^*]_{i:} - [\hat{\mathbf{A}}_{\tau}]_{i:}\|_2^2}_{\text{neighbor alignment regularizer}} \quad \Phi_{\tau} = \Phi - \alpha \nabla_{\Phi} \mathcal{L}(\mathcal{S}_{\tau}, f_{\theta, \Phi})$$

classification loss   neighbor alignment regularizer

to **separately capture** the generic knowledge shared across different tasks and those property-aware

# Algorithm: Selective update

---

**Algorithm 1** Meta-training procedure for PAR.

---

- 1: initialize  $\theta = \{\mathbf{W}_g, \mathbf{W}_a, \mathbf{W}_r\}$  and  $\Phi = \{\mathbf{W}_p, \mathbf{W}_c\}$  randomly; if a pretrained molecular encoder is available, take its parameter as  $\mathbf{W}_g$ ;
  - 2: **while** not done **do**
  - 3:   sample a batch of tasks  $\mathcal{T}_\tau$ ;
  - 4:   **for** all  $\mathcal{T}_\tau$  **do**
  - 5:     sample support set  $\mathcal{S}_\tau$  and query set  $\mathcal{Q}_\tau$  from  $\mathcal{T}_\tau$ ;
  - 6:     obtain molecular embedding  $\mathbf{g}_{\tau,i}$  for each  $\mathbf{x}_{\tau,i}$  by a graph-based molecular encoder;
  - 7:     adapt  $\mathbf{g}_{\tau,i}$  to be property-aware  $\mathbf{p}_{\tau,i}$  by (5);
  - 8:     initialize node embeddings as  $\mathbf{h}_{\tau,i}^{(0)} = \mathbf{p}_{\tau,i}$ ;
  - 9:     **for**  $t = 1, \dots, T$  **do**
  - 10:       estimate adjacency matrix  $\mathbf{A}_\tau^{(t)}$  of relation graph among molecules using  $\mathbf{h}_{\tau,i}^{(t-1)}$  by (6);
  - 11:       refine  $\mathbf{h}_{\tau,i}^{(t)}$  on the updated relation graph  $\mathbf{A}_\tau^{(t)}$  by (8);
  - 12:     **end for**
  - 13:     obtain class prediction  $\hat{\mathbf{y}}_{\tau,i}$  using  $\mathbf{h}_{\tau,i} = \mathbf{h}_{\tau,i}^{(T)}$ ;
  - 14:     evaluate training loss  $\mathcal{L}(\mathcal{S}_\tau, f_{\theta, \Phi})$  on  $\mathcal{S}_\tau$ ;
  - 15:     fine-tune  $\Phi$  as  $\Phi_\tau$  by (11);
  - 16:     evaluate testing loss  $\mathcal{L}(\mathcal{Q}_\tau, f_{\theta, \Phi_\tau})$  on  $\mathcal{Q}_\tau$ ;
  - 17:   **end for**
  - 18:   update  $\theta$  and  $\Phi$  by (12);
  - 19: **end while**
- 

Property-aware  
Molecular Embedding

Graph  
Construction

Selective update

MAML-based

# Outline

- Background
- Existing works
- The proposed approach
  - Property-Aware Relation networks (PAR)
  - Experiments
- Summary



# Experiment Setup

- Two sets of baselines
  - Methods with **graph-based encoder learned from scratch** including Siamese [Koch et al., 2015], ProtoNet [Snell et al., 2017], MAML [Finn et al., 2017], TPN [Liu et al., 2018], and EGNN [Kim et al., 2019], IterRefLSTM [Altae-Tran et al., 2017];
  - Methods which **leverage pretrained** graph-based molecular encoder including Pre-GNN [Hu et al., 2019], Meta-MGNN [Guo et al., 2021], and Pre-PAR which is our PAR equipped with Pre- GNN.

- Four datasets

Dataset	Tox21	SIDER	MUV	ToxCast
# Compounds	8014	1427	93127	8615
# Tasks	12	27	17	617
# Meta-Training Tasks	9	21	12	450
# Meta-Testing Tasks	3	6	5	167

- Link of our code: <https://github.com/tata1661/PAR-NeurIPS21>

# Overall Comparison

Method		Tox21		SIDER		MUV		ToxCast	
		10-shot	1-shot	10-shot	1-shot	10-shot	1-shot	10-shot	1-shot
Train directly	Siamese	80.40 <sub>(0.35)</sub>	65.00 <sub>(1.58)</sub>	71.10 <sub>(4.32)</sub>	51.43 <sub>(3.31)</sub>	59.96 <sub>(5.13)</sub>	50.00 <sub>(0.17)</sub>	-	-
	ProtoNet	74.98 <sub>(0.32)</sub>	65.58 <sub>(1.72)</sub>	64.54 <sub>(0.89)</sub>	57.50 <sub>(2.34)</sub>	65.88 <sub>(4.11)</sub>	58.31 <sub>(3.18)</sub>	63.70 <sub>(1.26)</sub>	56.36 <sub>(1.54)</sub>
	MAML	80.21 <sub>(0.24)</sub>	75.74 <sub>(0.48)</sub>	70.43 <sub>(0.76)</sub>	67.81 <sub>(1.12)</sub>	63.90 <sub>(2.28)</sub>	60.51 <sub>(3.12)</sub>	66.79 <sub>(0.85)</sub>	65.97 <sub>(5.04)</sub>
	TPN	76.05 <sub>(0.24)</sub>	60.16 <sub>(1.18)</sub>	67.84 <sub>(0.95)</sub>	62.90 <sub>(1.38)</sub>	65.22 <sub>(5.82)</sub>	50.00 <sub>(0.51)</sub>	62.74 <sub>(1.45)</sub>	50.01 <sub>(0.05)</sub>
	EGNN	81.21 <sub>(0.16)</sub>	79.44 <sub>(0.22)</sub>	72.87 <sub>(0.73)</sub>	70.79 <sub>(0.95)</sub>	65.20 <sub>(2.08)</sub>	62.18 <sub>(1.76)</sub>	63.65 <sub>(1.57)</sub>	61.02 <sub>(1.94)</sub>
	IterRefLSTM	81.10 <sub>(0.17)</sub>	80.97 <sub>(0.10)</sub>	69.63 <sub>(0.31)</sub>	71.73 <sub>(0.14)</sub>	49.56 <sub>(5.12)</sub>	48.54 <sub>(3.12)</sub>	-	-
Pre-train GNN	PAR	82.06 <sub>(0.12)</sub>	80.46 <sub>(0.13)</sub>	74.68 <sub>(0.31)</sub>	71.87 <sub>(0.48)</sub>	66.48 <sub>(2.12)</sub>	64.12 <sub>(1.18)</sub>	69.72 <sub>(1.63)</sub>	67.28 <sub>(2.90)</sub>
	Pre-GNN	82.14 <sub>(0.08)</sub>	81.68 <sub>(0.09)</sub>	73.96 <sub>(0.08)</sub>	73.24 <sub>(0.12)</sub>	67.14 <sub>(1.58)</sub>	64.51 <sub>(1.45)</sub>	73.68 <sub>(0.74)</sub>	72.90 <sub>(0.84)</sub>
	Meta-MGNN	82.97 <sub>(0.10)</sub>	82.13 <sub>(0.13)</sub>	75.43 <sub>(0.21)</sub>	73.36 <sub>(0.32)</sub>	68.99 <sub>(1.84)</sub>	65.54 <sub>(2.13)</sub>	-	-
	Pre-PAR	84.93 <sub>(0.11)</sub>	83.01 <sub>(0.09)</sub>	78.08 <sub>(0.16)</sub>	74.46 <sub>(0.29)</sub>	69.96 <sub>(1.37)</sub>	66.94 <sub>(1.12)</sub>	75.12 <sub>(0.84)</sub>	73.63 <sub>(1.00)</sub>

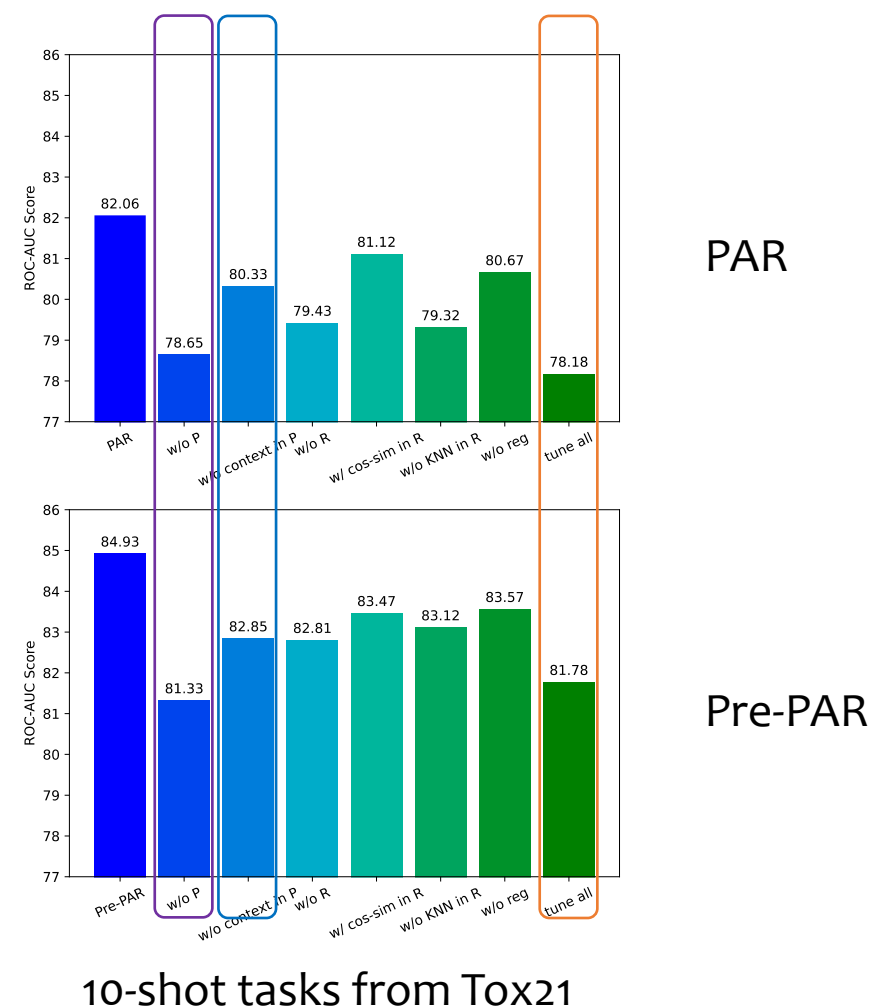
- Pre-PAR **consistently** obtains the **best** performance
- PAR outperforms among methods without pretrained GNNs

# Ablation Study

We further compare with

- **w/o P:** w/o **property-aware embedding**
- **w/o context in P:** w/o context  $b_{\tau,i}$  in P
- **w/o R:** w/o **adaptive relation graph learning**
- **w/ cos-sim in R:** use cosine similarity to obtain the adjacency matrix
- **w/o KNN in R:** w/o reducing the learned relation graph to KNN graph
- **w/o reg:** w/o the neighbor alignment regularizer
- **tune all:** **without selective update**

All components **are vital** to the success of PAR



# Case Study on 10 Molecules

Can PAR obtain **different** property-aware molecular **embeddings** and **relation graphs** for tasks containing **overlapping molecules** but evaluating **different properties**?

Table 5: The 10 molecules sampled from Tox21 dataset, which coexist in the three meta-testing tasks (the 10th task for SR-HSE, the 11th task for SR-MMP, and the 12th task for SR-p53).

Molecule		Label		
ID	SMILES	SR-HSE	SR-MMP	SR-p53
Mol-1	<chem>Cc1cccc(/N=N/c2ccc(N(C)C)cc2)c1</chem>	0	1	0
Mol-2	<chem>O=C(c1ccccc1)C1CCC1</chem>	1	0	0
Mol-3	<chem>C=C(C)[C@H]1CN[C@H](C(=O)O)[C@H]1CC(=O)O</chem>	0	0	1
Mol-4	<chem>c1ccc2sc(SNC3CCCCC3)nc2c1</chem>	1	1	0
Mol-5	<chem>C=CCSSCC=C</chem>	0	0	1
Mol-6	<chem>CC(C)(C)c1cccc(C(C)(C)C)c1O</chem>	0	1	0
Mol-7	<chem>C[C@@H]1CC2(OC3C[C@@]4(C)C5=CC[C@H]6C(C)(C)C(O[C@@H]7OC[C@@H](O)[C@H](O)[C@H]7O)CC[C@]67C[C@@]57CC[C@]4(C)C31)OC(O)C1(C)OC21</chem>	0	1	0
Mol-8	<chem>O=C(CCCCCC(=O)Nc1ccccc1)NO</chem>	0	0	1
Mol-9	<chem>CC/C=C\C/C=C\C/C=C\C\CCCCCCCC(=O)O</chem>	1	0	0
Mol-10	<chem>Cl[Si](Cl)(c1ccccc1)c1ccccc1</chem>	0	1	0

a fixed group of 10 molecules coexist in different meta-testing tasks

# Visualization

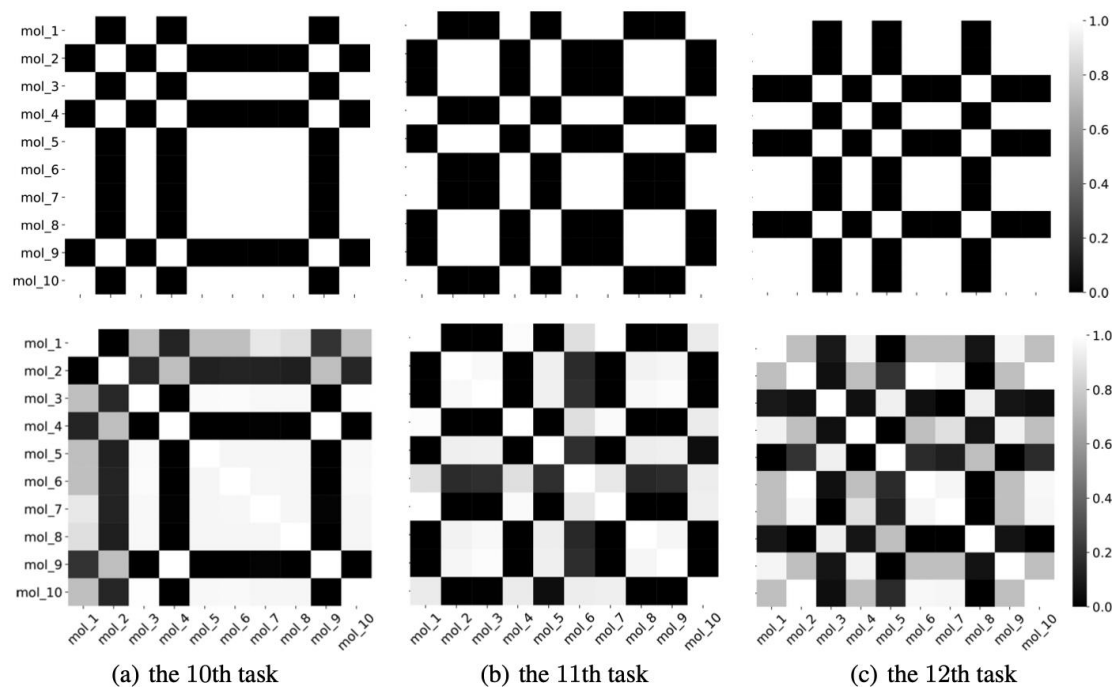


Figure 5: Comparison between  $\mathbf{A}_\tau^*$  computed using ground-truth labels (the first row) and adjacency matrix  $\mathbf{A}_\tau$  returned by PAR (the second row) for the ten molecules. We set  $[\mathbf{A}_\tau^*]_{ij} = 1$  if molecules  $\mathbf{x}_{\tau,i}$  and  $\mathbf{x}_{\tau,j}$  have the same label and 0 otherwise.

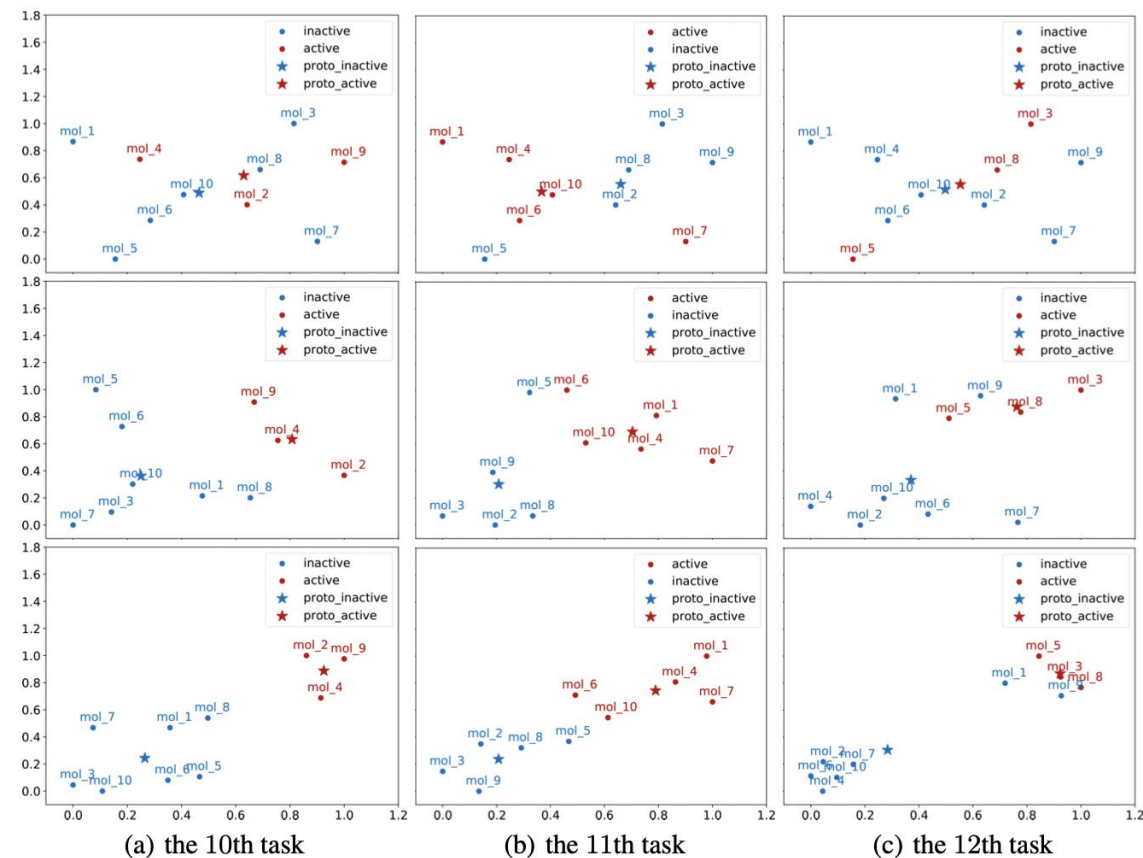


Figure 6: t-SNE visualization of  $\mathbf{g}_{\tau,i}$  (the first row),  $\mathbf{p}_{\tau,i}$  (the second row), and  $\mathbf{h}_{\tau,i}$  (the third row) of the ten molecules. Proto\_active (proto\_inactive) denotes the class prototype of active (inactive) class.

PAR can model relation graphs (left) and property-aware molecular embeddings (right)

# Outline

- Background
- Existing works
- The proposed approach
- Summary



# Summary

- Molecular property prediction (MPP) is a few-shot learning problem on graphs
- Existing works ignore molecule's difference in different tasks and molecule's relationship inside one task
- Our Property-Aware Relation networks (PAR) solves above problems with
  - Property-aware molecular embedding / Relation graph learning / Selective update
- PAR significantly outperforms existing works

# FSL: Learning resources

## A CSUR Survey

- Generalizing from a Few Examples: A Survey on Few-Shot Learning. ACM CSUR. 2020

- Show you a roadmap


## A GitHub Repo

- <https://github.com/tata1661/FSL-Mate>
- Update-to-date research papers

## A Toolbox

- Support various applications and platforms
- Built with Baidu

SURVEY  
June 2020



Generalizing from a Few Examples: A Survey on Few-shot Learning  
Yaqing Wang, Quanming Yao, James T. Kwok, Lionel M. Ni  
ACM Computing Surveys (CSUR), Volume 53, Issue 3 • May 2021, Article No.: 63, pp 1–34 • <https://doi.org/10.1145/3386252>  
Machine learning has been highly successful in data-intensive applications but is often hampered when the data set is small. Recently, Few-shot Learning (FSL) is proposed to tackle this problem. Using prior knowledge, FSL can rapidly generalize to new ...  
314 9,508  
Get Access

tata1661 / FSL-Mate Public
Unwatch 44 Fork 254 Starred 1.3k
Code Issues 6 Pull requests 3 Actions Projects Wiki Security Insights
master 2 branches 0 tags Go to file Add file Code About
tata1661 Add PaddleFSL paper and bib, fix ... 676f5a3 27 days ago 60 commits
FewShotPapers PaddleFSL updates to v1.1.0 (add support for ... last month
PaddleFSL Add PaddleFSL paper and bib, fix bugs in exa... 27 days ago
FSL-Mate: A collection of resources for few-shot learning (FSL).
deep-learning paper papers one-shot-learning paddlepaddle meta-learning few-shot

		learn2learn	Torchmeta	keras-fsl	mmfewshot	Libfewshot	PaddleFSL
CV	MAML	✓	✓	✗	✓	✓	✓
	ProtoNet	✓	✓	✗	✓	✓	✓
	RelationNet	✗	✗	✗	✓	✓	✓
Appli- cation NLP	Siamese	✗	✗	✓	✗	✗	✓
	PET	✗	✗	✗	✗	✗	✓
	P-Tuning	✗	✗	✗	✗	✗	✓
BIO	MAML	✗	✗	✗	✗	✗	✓
	PAR	✗	✗	✗	✗	✗	✓
Unit Test		✓	✓	✗	✓	✗	✓
Document		✓	✓	✓	✓	✓	✓
Platform	CPU	✓	✓	✓	✓	✓	✓
	GPU	✓	✓	✓	✓	✓	✓
	Cluster	✗	✗	✗	✗	✗	✓
Activity	Download	45182	131963	5484	913	N/A	24441
	Star	1725	1576	172	359	496	1247
	Fork	249	203	29	42	87	240

Table 1: Comparison of PaddleFSL with other popular FSL toolboxes. The activity statistics is collected from the respective GitHub pages on 2022/3/18.



# Future Works

- Transfer learning across datasets
- Few-shot molecule regression
- Design GNNs for better molecular representations

Thanks! Questions?