# Modeling Chemical Reactions

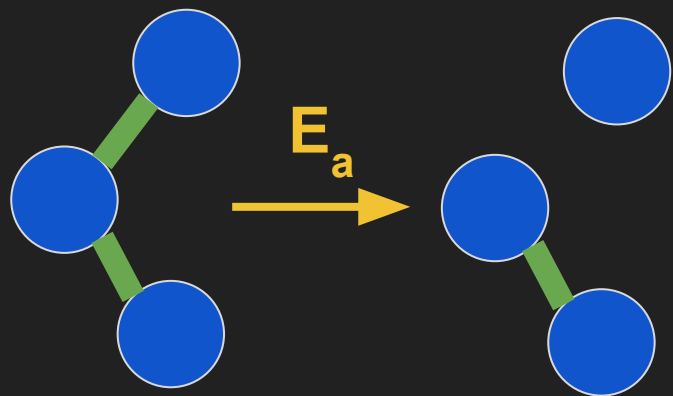**Intern:** Shaun Tan *(me)*

**Supervisors:**
Chieu Hai Leong, Chong Yihui, Alvin Liew

# About me

- NUS, year 2 Computer Science student
- Came from Diploma of Applied Chemistry

- Started internship with only introductory knowledge in neural networks
- Interested in doing **SWE** and/or **ML**

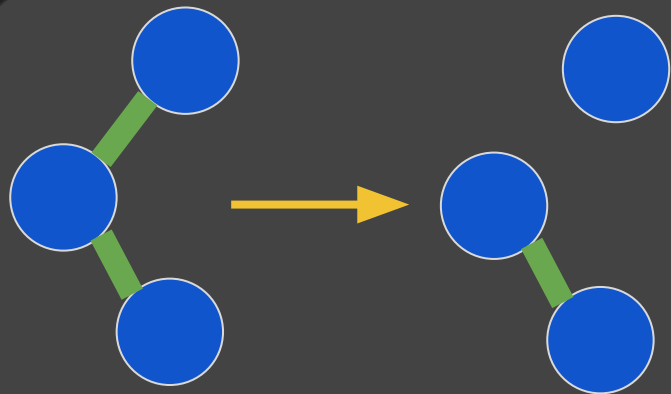# Objective



**Reactant**                    **Product**

Predicting **activation energy** from the reaction's **molecules** *(ie. reactants & products)*

# Objective



**Reaction properties**
- Multiple molecules
- Has reactants & products

- Activation energy
- Rate of reaction
- Yield %

VS

- Water solubility
- Toxicity

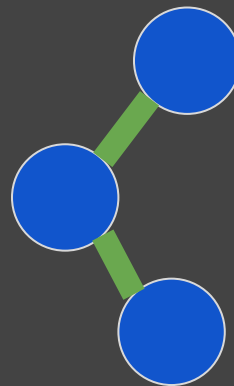**Molecular properties**
- 1 molecule only

# Objective

Can we <u>train on **molecular**-properties</u>, then <u>finetune on **reaction**-properties</u>?

- **Reason:** there's a lot more molecular datasets than reaction datasets

Can a model for that's good at predicting <u>molecular-property</u> be **adapted** to predict <u>reaction properties</u>?
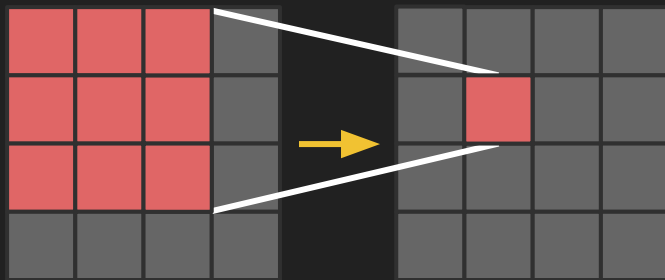


- Water solubility
- Toxicity

<u>**Molecular properties**</u>

- 1 molecule only

# **Context:** Graph Neural Networks (GNN)
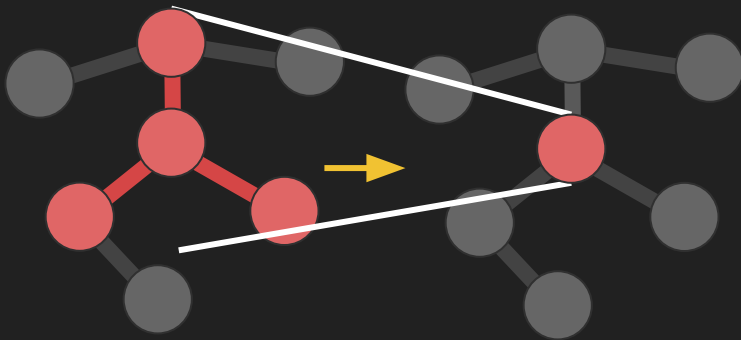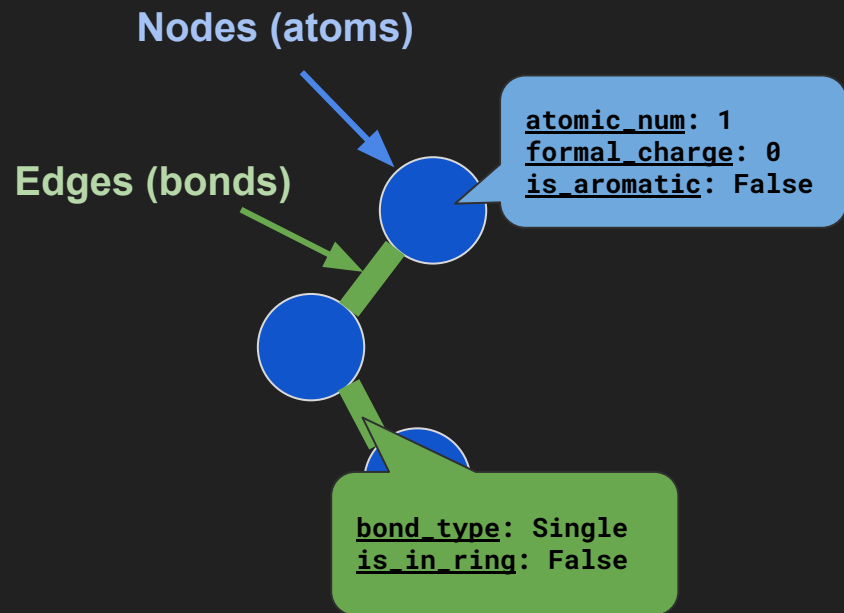


**CNN**

Like how **CNN** aggregate adjacent pixels in an image *(ie. convolving)*

**GNN**

**GNN** aggregate adjacent nodes & edges in a graph *(ie. message-passing)*

# Context: GNN for molecules

**Nodes (atoms)**

**Edges (bonds)**

```
atomic_num: 1
formal_charge: 0
is_aromatic: False
```

```
bond_type: Single
is_in_ring: False
```

Molecules can modelled as graphs
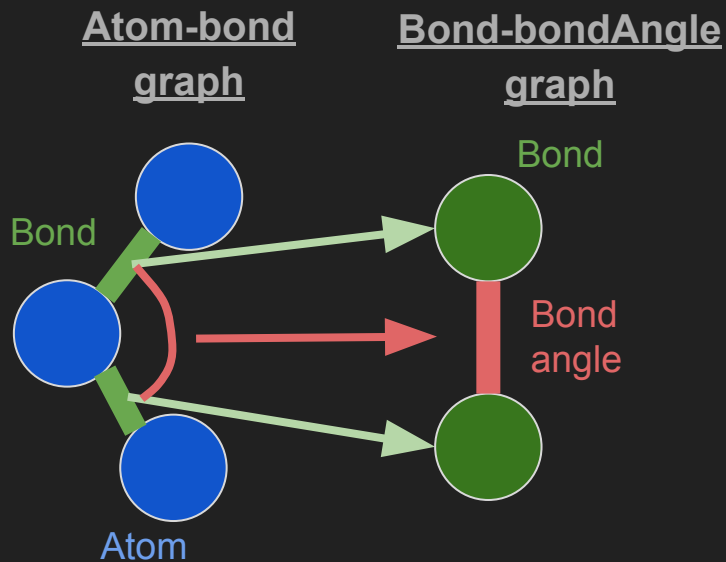
# Using SOTA architectures

- We used **GeoGNN (aka GEM)**, a State-Of-The-Art (SOTA) architecture for predicting **molecular**-properties

**Table 1 | Overall performance for regression tasks and classification tasks**
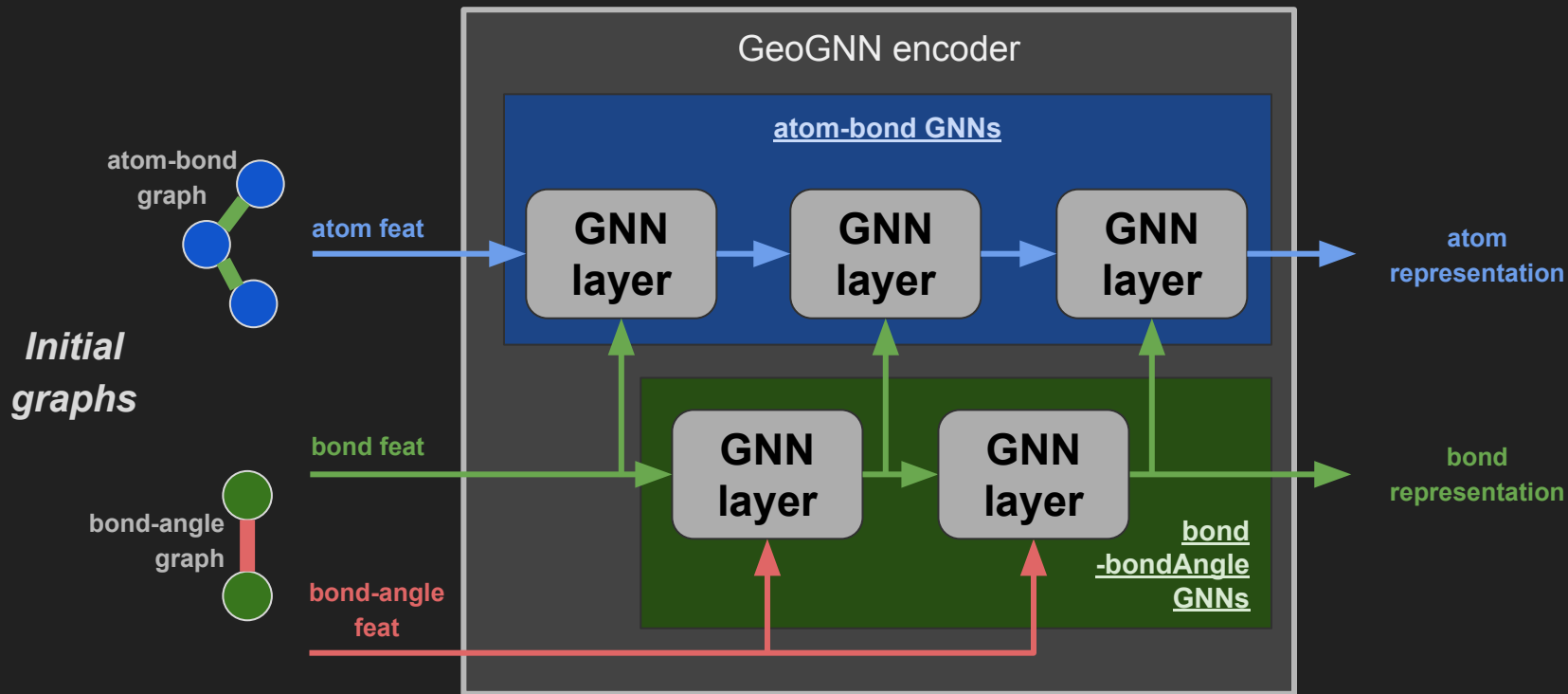
**Regression (lower is better)**

| | RMSE | | | | MAE | |
|---|---|---|---|---|---|---|
| Dataset | ESOL | FreeSolv | Lipo | QM7 | QM8 | QM9 |
| No. molecules | 1,128 | 642 | 4,200 | 6,830 | 21,786 | 133,885 |
| No. prediction tasks | 1 | 1 | 1 | 1 | 12 | 12 |
| D-MPNN [43] | $1.050_{(0.008)}$ | $2.082_{(0.082)}$ | $^a0.683_{(0.016)}$ | $103.5_{(8.6)}$ | $0.0190_{(0.0001)}$ | $0.00814_{(0.00001)}$ |
| AttentiveFP [44] | $^a0.877_{(0.029)}$ | $^a2.073_{(0.183)}$ | $0.721_{(0.001)}$ | $^a72.0_{(2.7)}$ | $^a0.0179_{(0.0001)}$ | $^a0.00812_{(0.00001)}$ |
| N-Gram$_{RF}$ [45] | $1.074_{(0.107)}$ | $2.688_{(0.085)}$ | $0.812_{(0.028)}$ | $92.8_{(4.0)}$ | $0.0236_{(0.0006)}$ | $0.01037_{(0.00016)}$ |
| N-Gram$_{XGB}$ [45] | $1.083_{(0.082)}$ | $5.061_{(0.744)}$ | $2.072_{(0.030)}$ | $81.9_{(1.9)}$ | $0.0215_{(0.0005)}$ | $0.00964_{(0.00031)}$ |
| PretrainGNN [11] | $1.100_{(0.006)}$ | $2.764_{(0.002)}$ | $0.739_{(0.003)}$ | $113.2_{(0.6)}$ | $0.0200_{(0.0001)}$ | $0.00922_{(0.00004)}$ |
| GROVER$_{base}$ [4] | $0.983_{(0.090)}$ | $2.176_{(0.052)}$ | $0.817_{(0.008)}$ | $94.5_{(3.8)}$ | $0.0218_{(0.0004)}$ | $0.00984_{(0.00055)}$ |
| GROVER$_{large}$ [4] | $0.895_{(0.017)}$ | $2.272_{(0.051)}$ | $0.823_{(0.010)}$ | $92.0_{(0.9)}$ | $0.0224_{(0.0003)}$ | $0.00986_{(0.00025)}$ |
| **(best)** GEM | $\mathbf{0.798}_{(0.029)}$ | $\mathbf{1.877}_{(0.094)}$ | $\mathbf{0.660}_{(0.008)}$ | $\mathbf{58.9}_{(0.8)}$ | $\mathbf{0.0171}_{(0.0001)}$ | $\mathbf{0.00746}_{(0.00001)}$ |

# Using SOTA architectures

**Atom-bond graph**

**Bond-bondAngle graph**

Bond

Bond

Atom

Bond

Bond angle

They proposed using 2 graphs to represent a molecule

# GeoGNN architecture

# GeoGNN architecture

# Adapting GeoGNN for reactions

# Results

# Performance

# Performance



Training / Validation losses during training

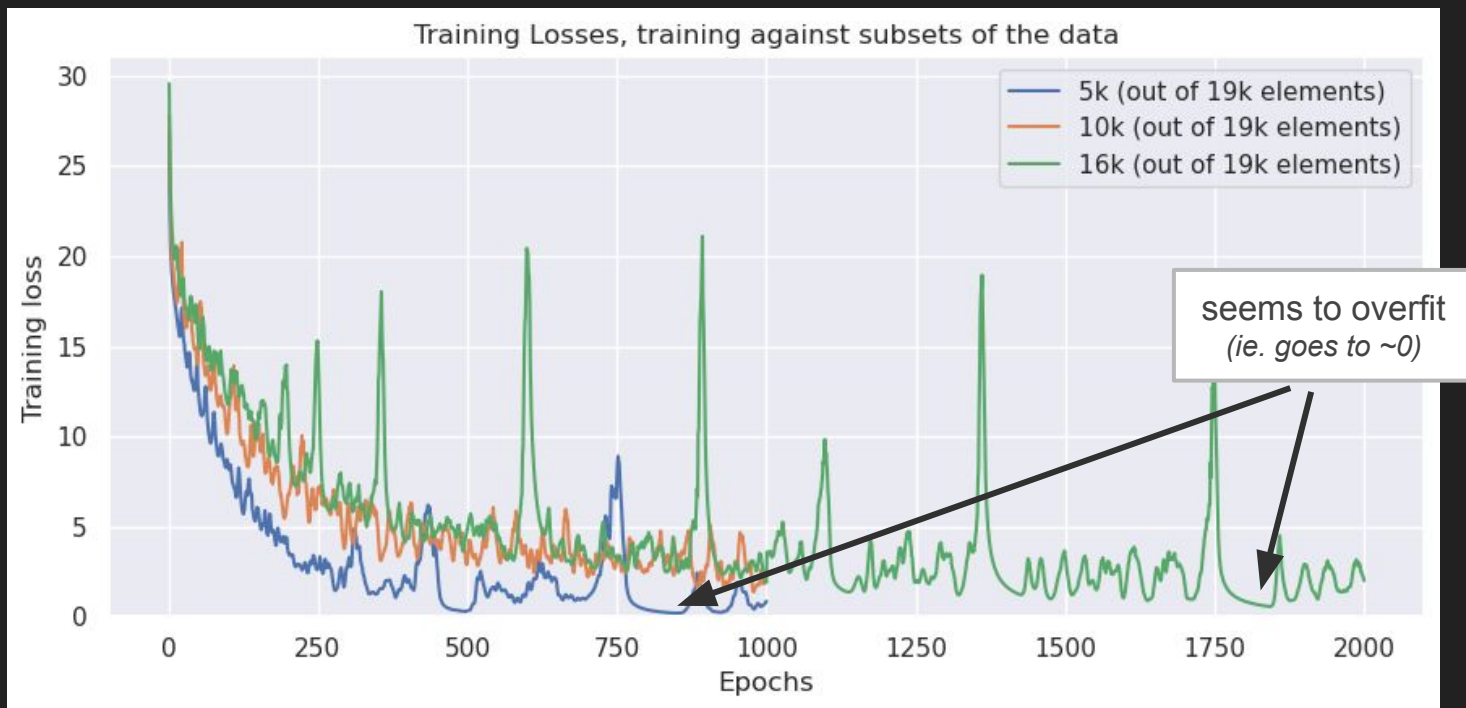# **Debugging:** Checking model's capacity

- **Intentionally** try to overfit  *(ie. converges to ~0 loss)*
  **smaller** dataset sizes  *(eg. 5, 10, 16k out of the 19k dataset)*

- If **can** overfit        =   model has the capacity to represent the dataset
- If **cannot** overfit   =   model is not powerful enough
                                        *(or I coded something wrongly)*

# **Debugging:** Checking model's capacity



Training Losses, training against subsets of the data

- 5k (out of 19k elements)
- 10k (out of 19k elements)
- 16k (out of 19k elements)

seems to overfit
*(ie. goes to ~0)*

**Thus:**
model seems to have the capacity to represent the dataset

# Tweaking the architecture



Training / Validation losses, with tweaked architectures

**validation** loss

**training** loss

Legend:
- Original
- with attention-pooling
- with attention-GNN (GAT)

Doesn't seem to improve performance

# Next steps

# **Next steps:** Tweaking combing/pooling



**Reaction graphs**

concatenated reactant & product feats

GeoGNN encoder

reactant & product atom repr

Downstream model

split atom reprs

reactant

product

**Combine + Graph Pool**

MLP

**reaction** property prediction

Continue experimenting on how to best combine/pool the reactant + product atom-representations

# **Next steps:** Pretraining on molecular-properties

- Pretrain on molecular-property datasets
  then finetune on activation-energy reaction datasets
  to see if the pretraining helps

# Challenges
# I faced

# Challenges I faced

- Recoded GeoGNN's architecture in PyTorch + DGL:
  - becuz their original code used a lesser-known library – paddlepaddle
  - and their code was very messy, w/o documentations
    both of which made it hard to modify/adapt



  - recoded to gain experience in working with GNN / PyTorch

# Challenges I faced

- Learnt PyTorch and DGL libraries
- To keep code clean, coded with type hintings, docstrings, PyTorch-Lightning

```python
class SqrtGraphNorm(torch.nn.Module):
    """
    Applies graph normalization, where each node features is divided by
    sqrt(num_of_nodes) for each graph in batched graph created by `dgl.batch`.

    This is a PyTorch + DGL equivalent of GeoGNN's `GraphNorm`:
    https://github.com/PaddlePaddle/PaddleHelix/blob/e93c3e9/pahelix/networks/gnn_
    """

    def forward(self, batched_graph: DGLGraph, node_feats: Tensor) -> Tensor:
        """
        Args:
            batched_graph (DGLGraph): Batched (or unbatched) DGL graph created by
            node_feats (Tensor): The input node features.

        Returns:
            Tensor: The node features that's been normalized via dividing by sqrt(
        """
        batch_num_of_nodes = batched_graph.batch_num_nodes()
        assert isinstance(batch_num_of_nodes, Tensor)
```

```python
class GeoGNNLightningModule(ABC, pl.LightningModule):

    def training_step(self, batch: GeoGNNBatch, batch_idx: int) -> Tensor:
        atom_bond_batch_graph, bond_angle_batch_graph, labels = batch
        pred = self.forward(atom_bond_batch_graph, bond_angle_batch_graph)
        self._train_step_values.append((pred, labels))
        loss = self.loss_fn(pred, labels)

        # Log raw unstandardized MSE loss to the progress bar and logger.
        self.log("train_raw_std_mse_loss", loss, on_step=False, on_epoch=True,

        return loss

    def predict_step(self, batch: GeoGNNBatch | tuple[DGLGraph, DGLGraph], batc
        atom_bond_batch_graph, bond_angle_batch_graph, *_ = batch
        pred = self.forward(atom_bond_batch_graph, bond_angle_batch_graph)
        pred = self.scaler.inverse_transform(pred)
        return pred

    def validation_step(self, batch: GeoGNNBatch, batch_idx: int) -> Tensor:
        atom_bond_batch_graph, bond_angle_batch_graph, labels = batch
        pred = self.forward(atom_bond_batch_graph, bond_angle_batch_graph)
```
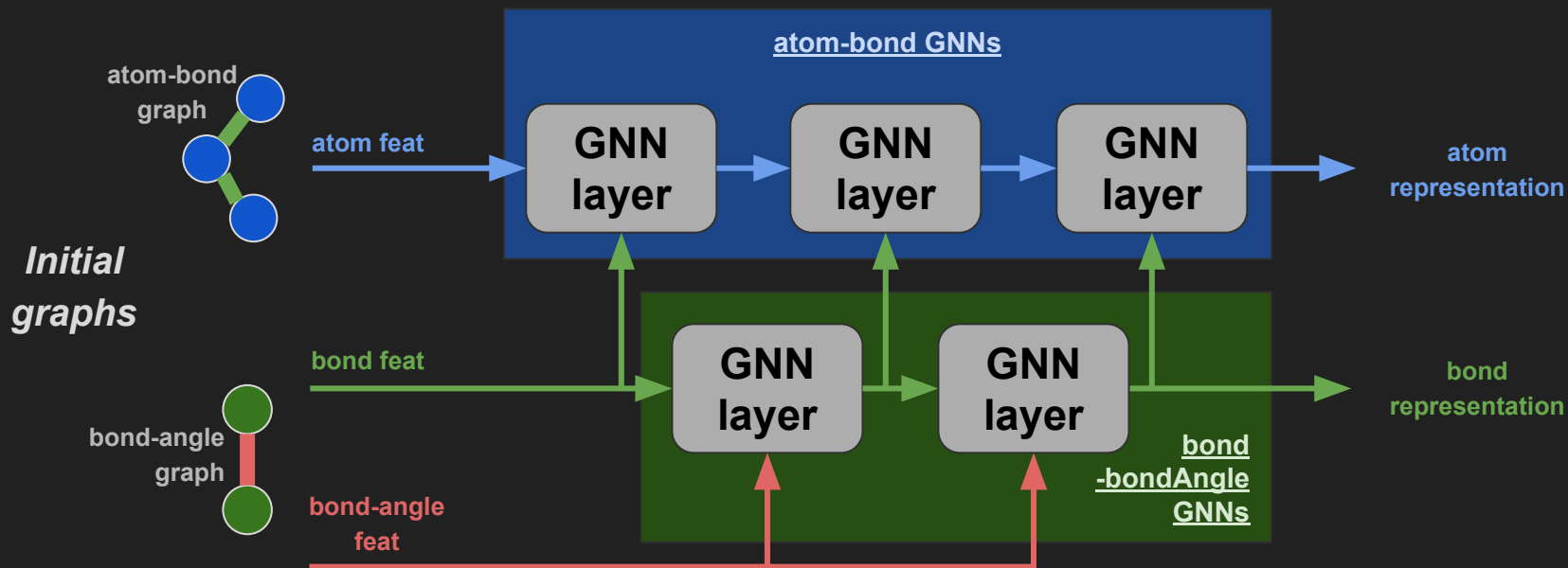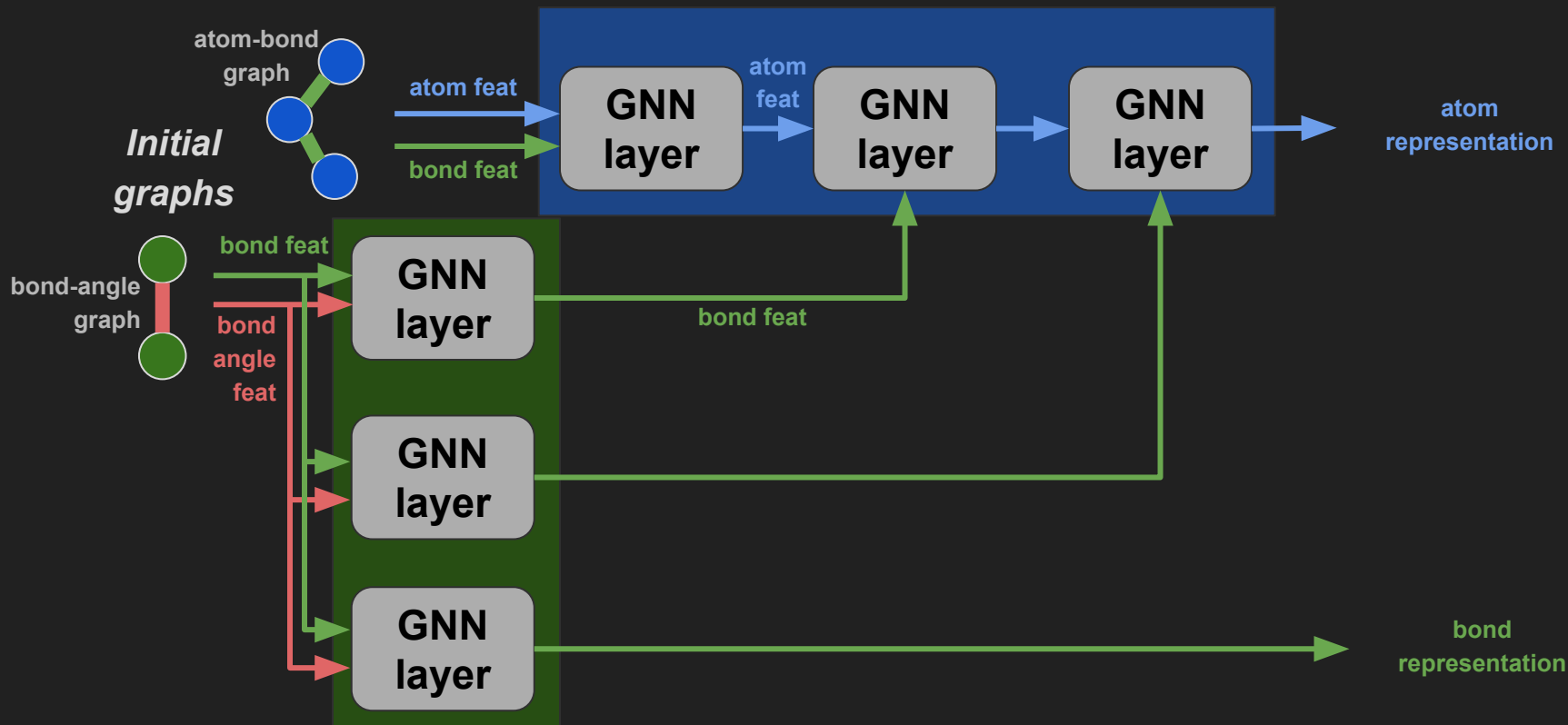
# Challenges I faced

- The model GeoGNN coded
  wasn't what they mentioned in the paper

# What they **SAID** they did:

# What they **ACTUALLY** did:

~ fin ~