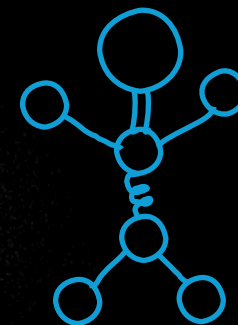# Learning a Molecule with GNNs

An interactive understanding of expressive GNNs for molecular representation, and how to scale them to infinity 🚀
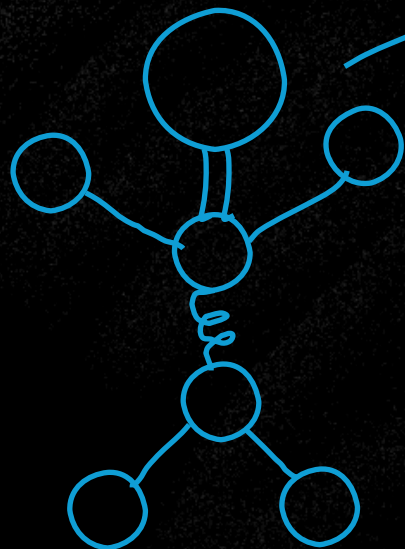
A presentation by Dominique Beaini

Research lead at Valence Labs

Adjunct Prof at Université de Montreal

Associate Industry member at Mila – Quebec AI institute

# Meeting Graphy 👋

Hello everyone 👋! I'm Dom's assistant for today!

Let's visit the molecular graph world together!

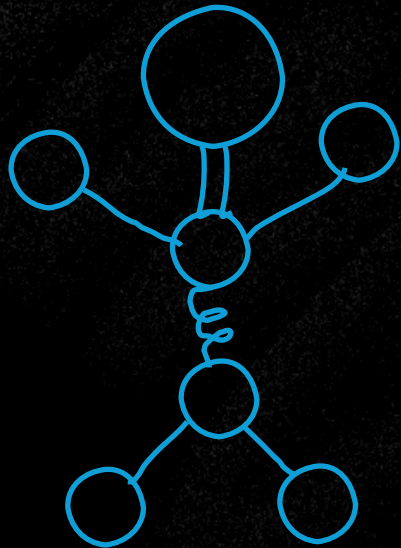We'll first learn what are graphs and how to manipulate them

Then we'll look into standard GNNs graph neural networks

And how to build more expressive GNNs for molecules

Then, we will scale a Graph Transformer together

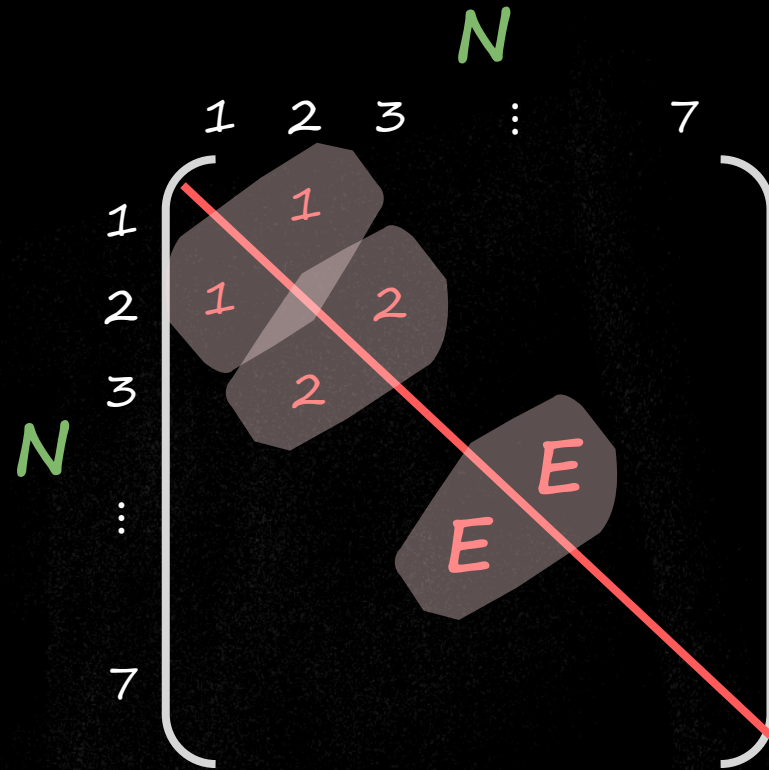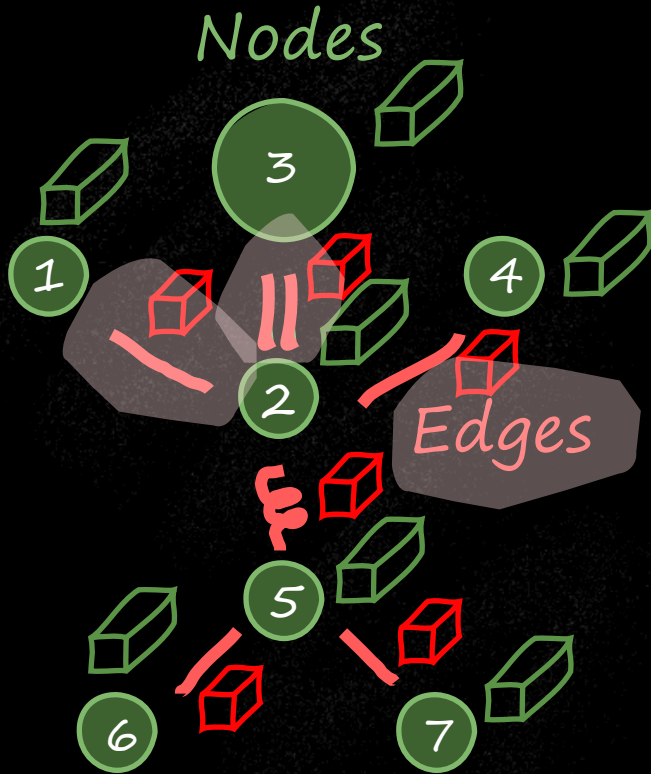Finally, we'll work through some applications

# Anatomy of Graphy 🦴

# Anatomy of Graphy 🦴

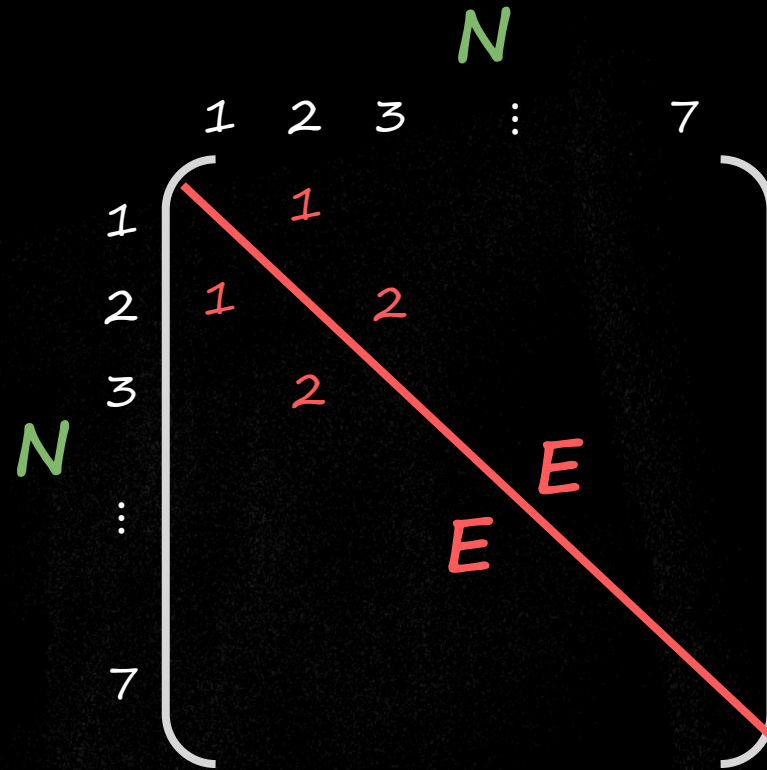Nodes

Edges

N

| | 1 | 2 | 3 | ⋮ | 7 |
|---|---|---|---|---|---|
| 1 | | 1 | | | |
| 2 | 1 | | 2 | | |
| 3 | | 2 | | | |
| ⋮ | | | | E | E |
| 7 | | | | | |

N

Adjacency matrix

# Anatomy of Graphy

Nodes

Edges

N

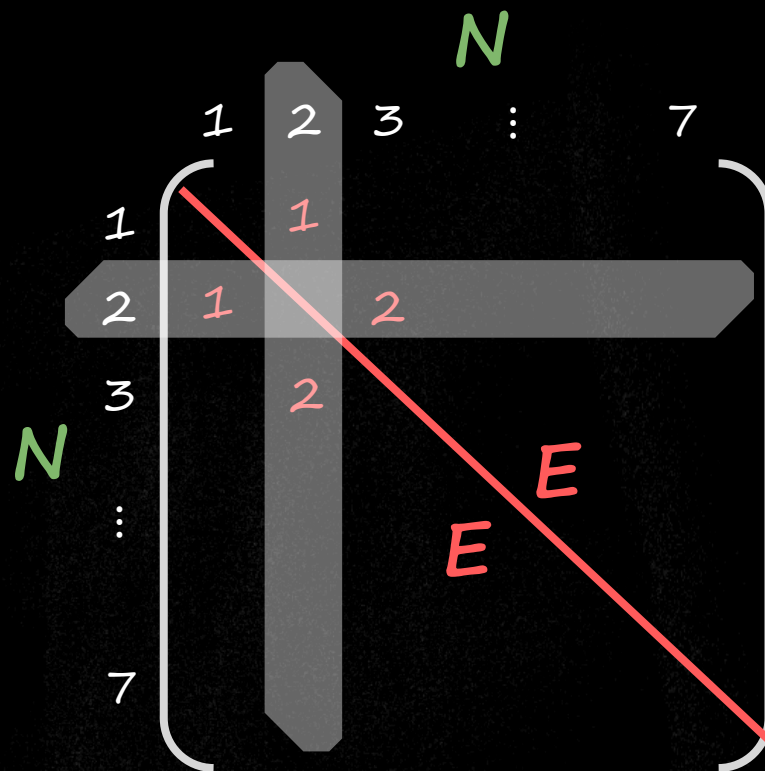| | 1 | 2 | 3 | ⋮ | 7 |
|---|---|---|---|---|---|
| 1 | | 1 | | | |
| 2 | 1 | | 2 | | |
| 3 | | 2 | | | |
| ⋮ | | | | E | |
| 7 | | | E | | |

N

Adjacency matrix

N

Node feature matrix

e

E

Edge feature matrix

# Permutation invariance 🔀

Nodes

Edges

$N$

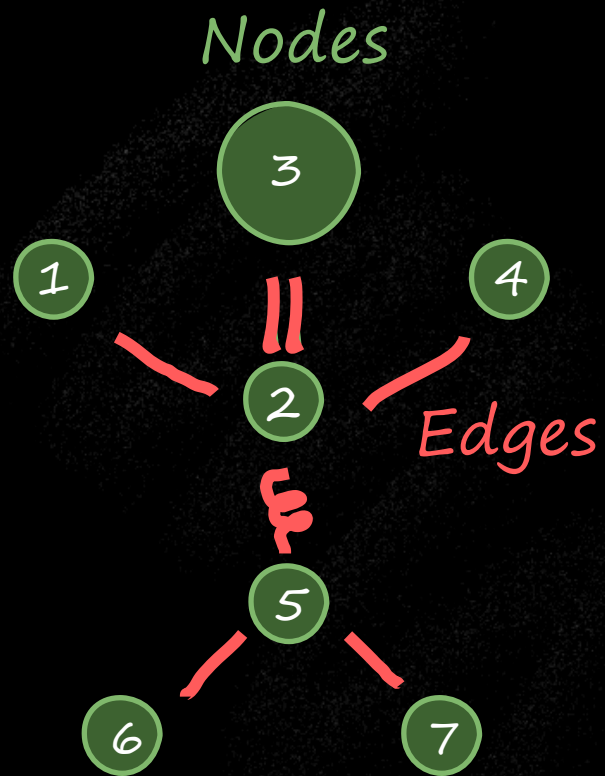|     | 1 | 2 | 3 | ⋮ | 7 |
|-----|---|---|---|---|---|
| 1   |   | 1 |   |   |   |
| 2   | 1 |   | 2 |   |   |
| 3   |   | 2 |   |   |   |
| ⋮   |   |   |   |   |   |
| 7   |   |   |   |   |   |

$N$

$E$

$E$

Adjacency matrix

$X$

$N$

Node feature matrix

$e$

$E$

Edge feature matrix

# Laplacian matrix 🧲

**Nodes**



**Edges**

$$N$$

| | 1 | 2 | 3 | : | 7 |
|---|---|---|---|---|---|
| 1 | | | *1* | | |
| 2 | *1* | | *2* | | |
| 3 | | *2* | | | |
| : | | | | *E* | |
| 7 | | | | | |

$N$

$$\Sigma \Longrightarrow$$

Degree Matrix

| 1 | | |
|---|---|---|
| | 3 | |
| | | 1 |

Adjacency matrix

**A**

**D**

$$L = D - A$$

# Anatomy of a molecule 💀 🦴

Nodes

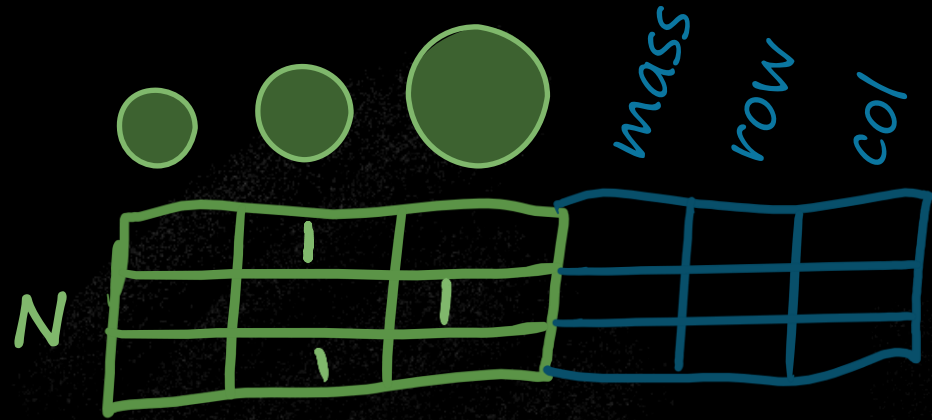Edges

Molecules are *defined* as graphs!

# Anatomy of atoms and bonds ⚛

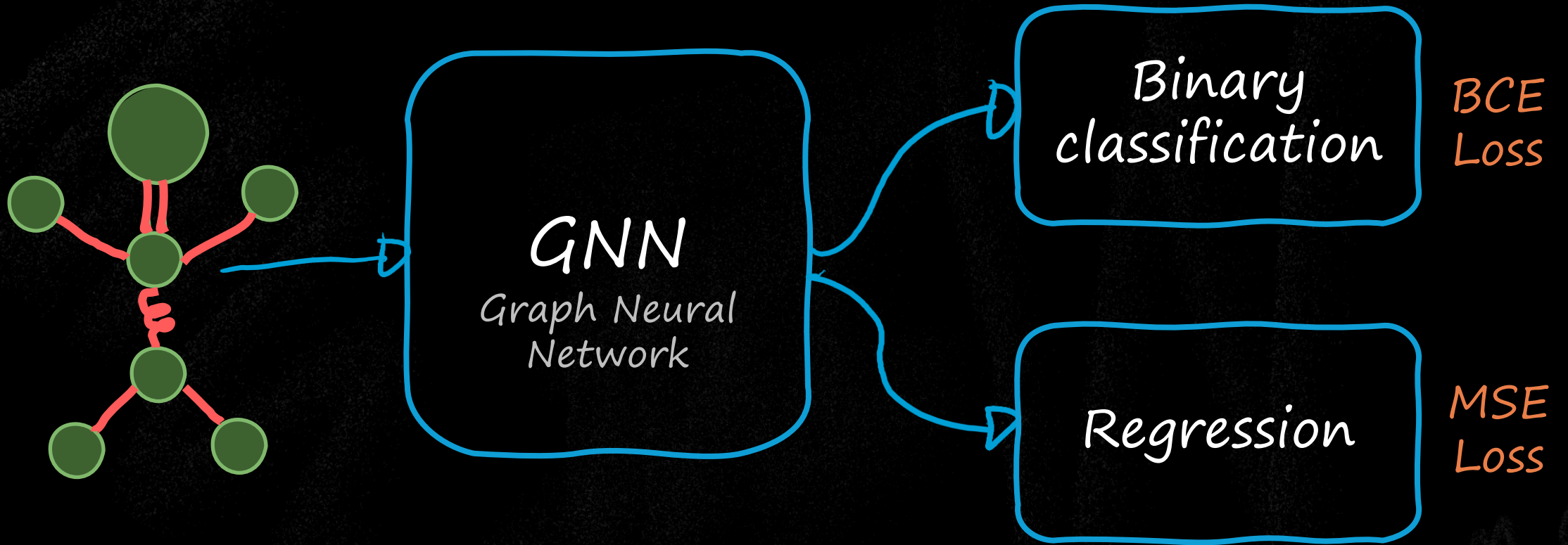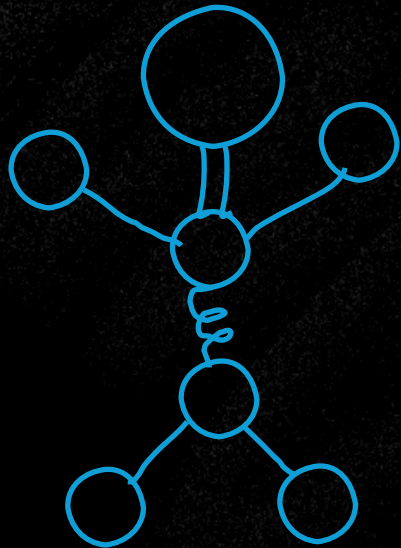N — Node feature matrix

mass · row · col

E — Edge feature matrix

length · chirality · rotatable

# Designing a GNN 💡

- How do we deal with permutation invariance?

- And the varying #nodes / neighbors?

- And the isotropy or lack of direction?
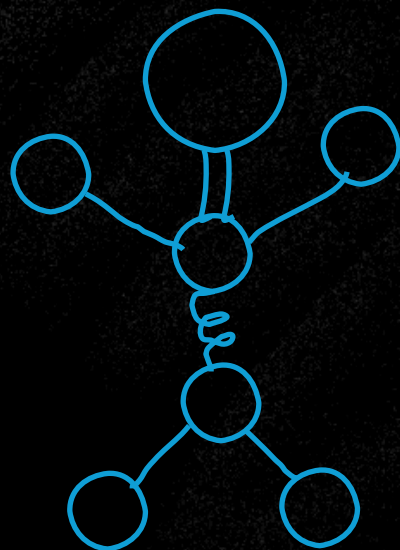
- And the expressivity?

# Designing a GNN 💡

## permutation invariance

- Apply MLPs on each node
- Pass the features on the edges

## varying #nodes

- Aggregation (mean/max/sum)
- Pooling (mean/max/sum)

## lack of direction

- Ignore that it is a problem 😬

## expressivity

- Make some weird proofs 🥴

*Don't worry, I'll show you a better way*

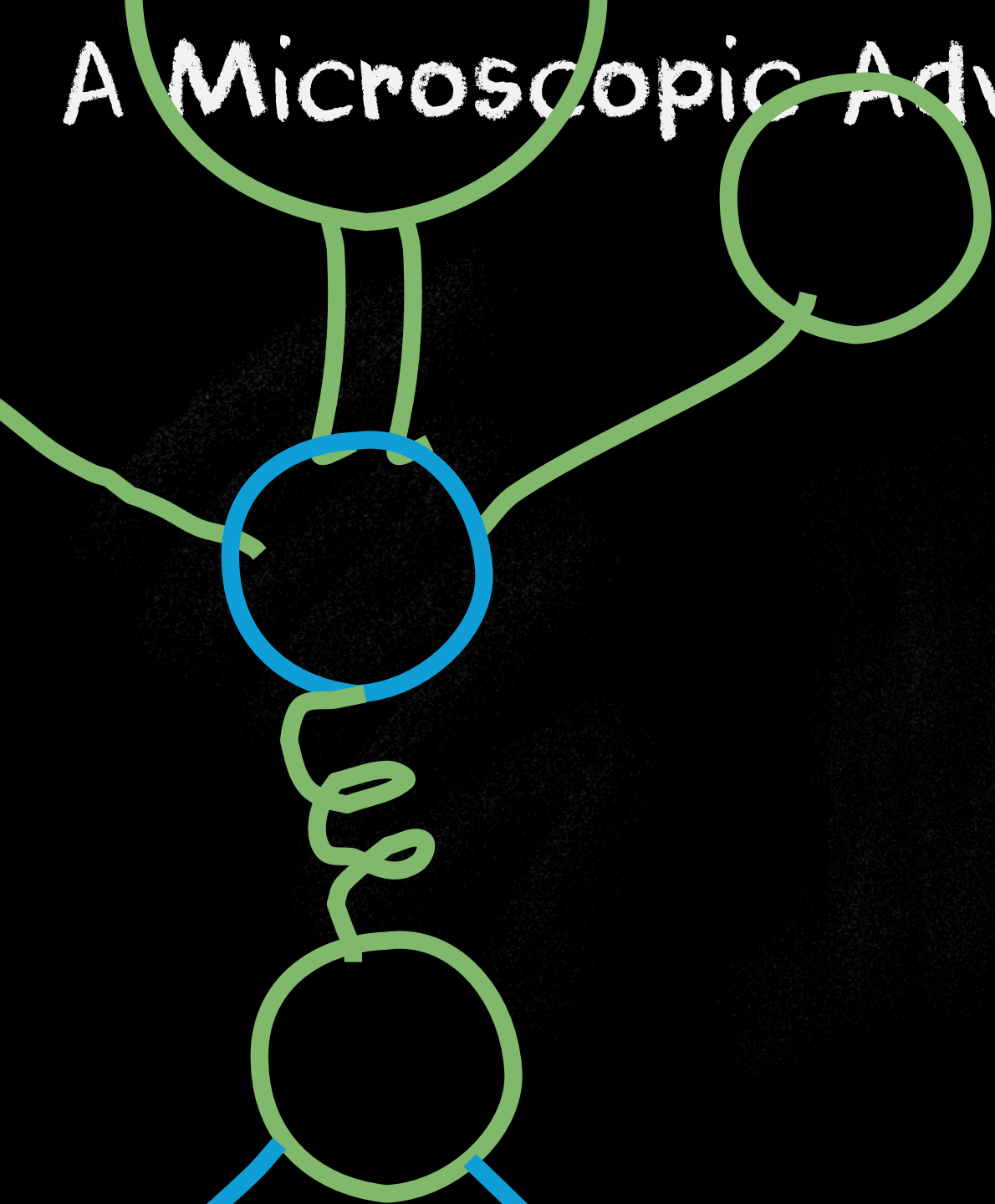# A Microscopic Adventure 🔬

Join me on a microscopic adventure!
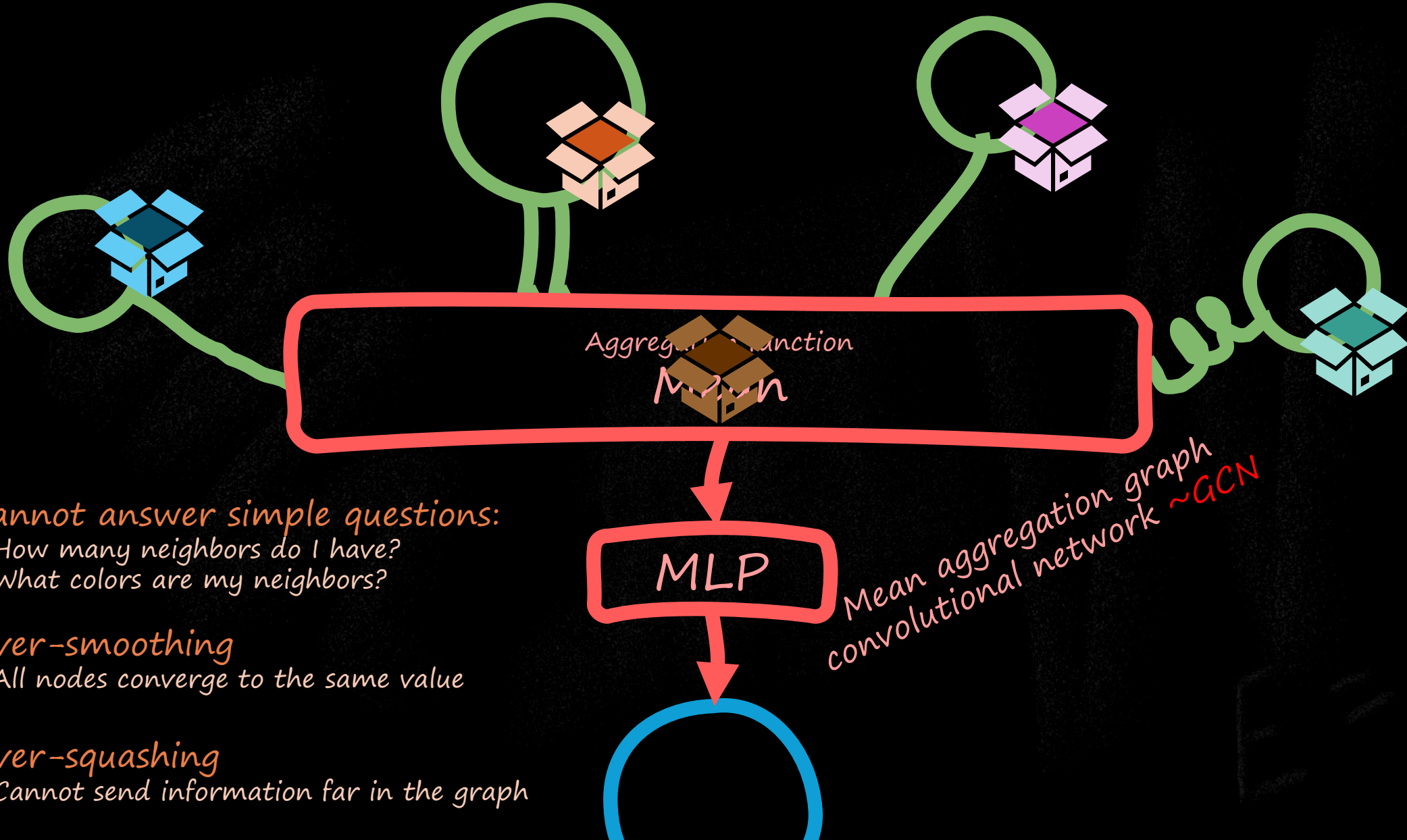Jump in the black hole to see the perspective of a node
Let's go, don't be afraid!

# A Microscopic Adventure 🔬

# Mean Aggregation Conv 💩

Aggregation function
Mean

Cannot answer simple questions:
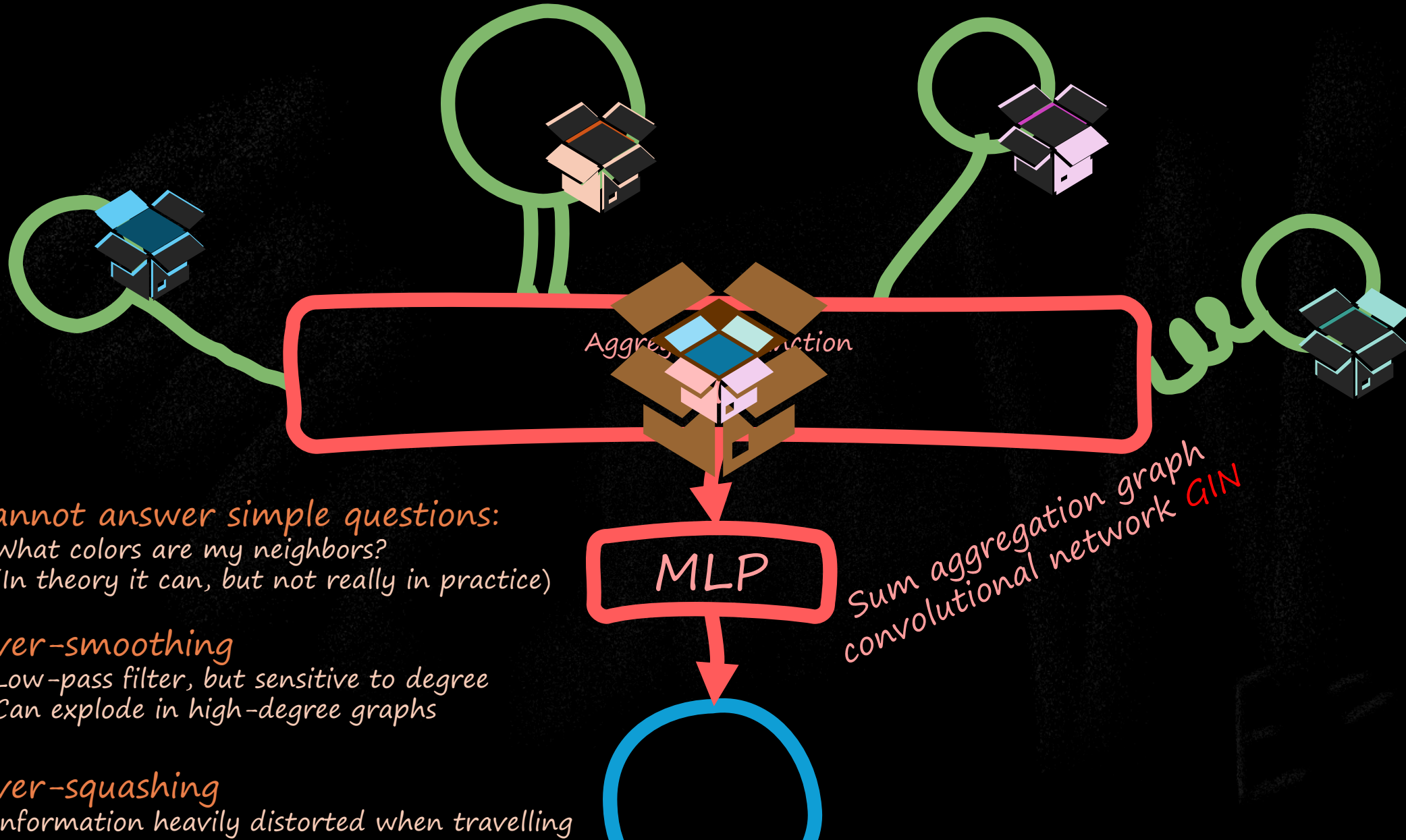How many neighbors do I have?
What colors are my neighbors?

Over-smoothing
All nodes converge to the same value

Over-squashing
Cannot send information far in the graph

MLP

Mean aggregation graph
convolutional network ~GCN

# Sum Aggregation Conv 🧨



Aggregation function

**Cannot answer simple questions:**
What colors are my neighbors?
(In theory it can, but not really in practice)

**Over-smoothing**
Low-pass filter, but sensitive to degree
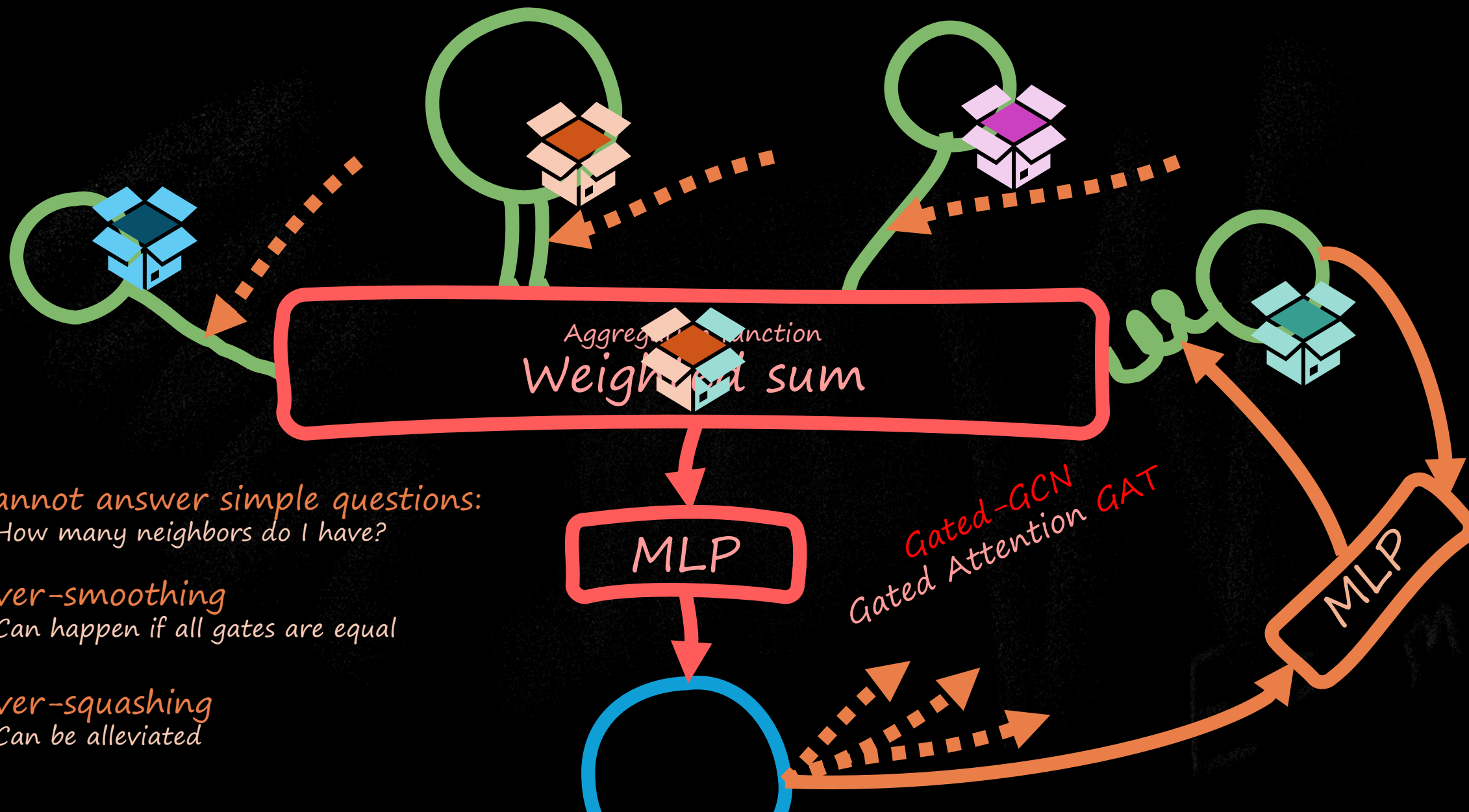Can explode in high-degree graphs

**Over-squashing**
Information heavily distorted when travelling

MLP

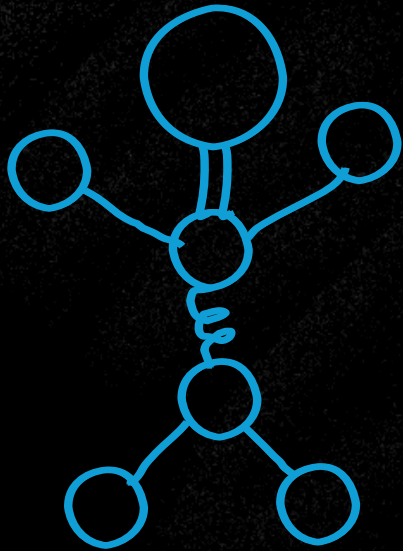Sum aggregation graph
convolutional network GIN

# Gating 🚪



Aggregation function
**Weighted sum**

MLP

Cannot answer simple questions:
How many neighbors do I have?

Over-smoothing
Can happen if all gates are equal

Over-squashing
Can be alleviated

Gated-GCN
Gated Attention GAT

MLP

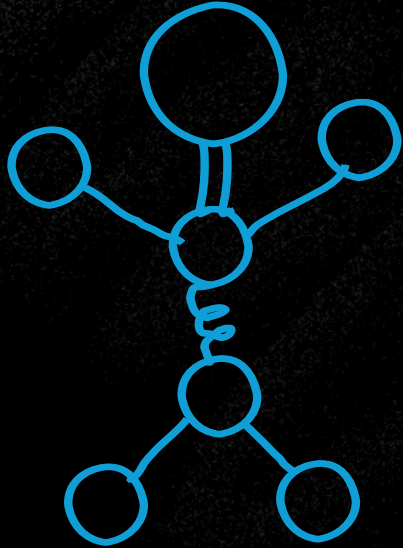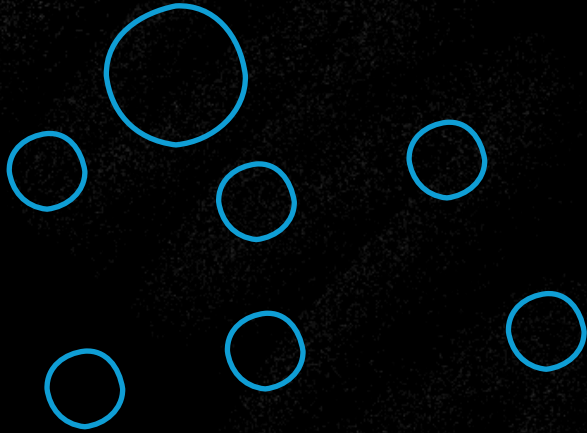# Attention is all you need
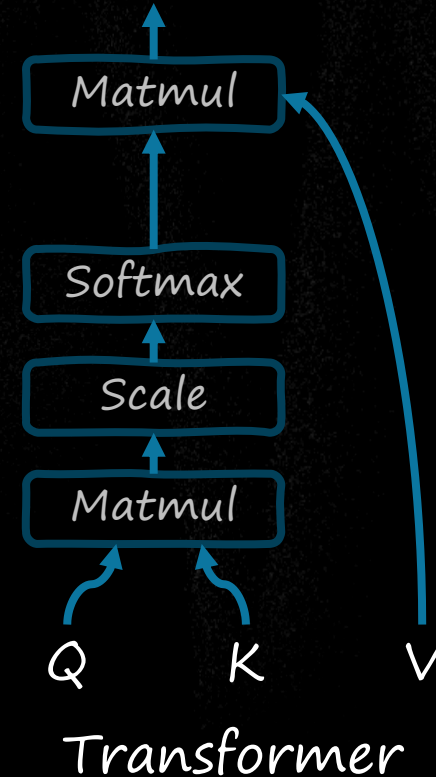
# Attention is all you need

- ## ATTENTION IS ALL YOU NEED
  - ### Certain conditions apply
    - #### Read the fineprints for more details
      - You also need good positional and structural encodings, ideally a biased attention, lots and lots of long-range data, ...
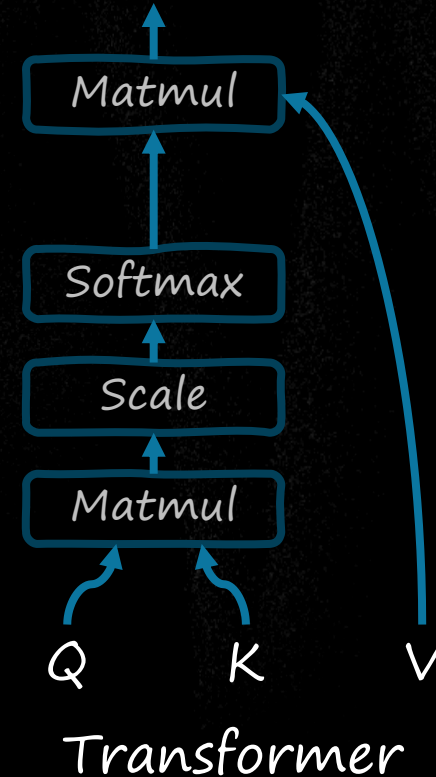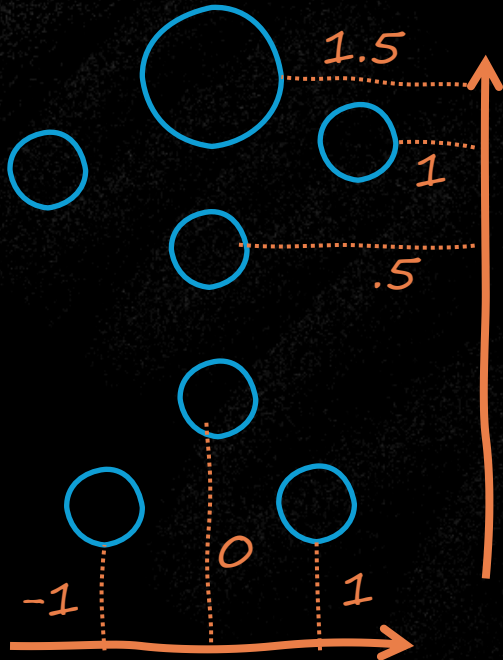
Matmul

Softmax

Scale

Matmul

Q          K          V

Transformer

# Attention is all you need

- **ATTENTION IS ALL YOU NEED**
  - *Certain conditions apply*
    - *Read the fineprints for more details*
      - You also need good positional and structural encodings, ideally a biased attention, lots and lots of long-range data, ...

What's happening to me!!!
I'm being permuted!

Matmul

Softmax

Scale

Matmul

Q      K      V

Transformer

# Fineprints of Attention – Position

- positional and structural encodings
- biased attention
- lots and lots of long-range data



Matmul

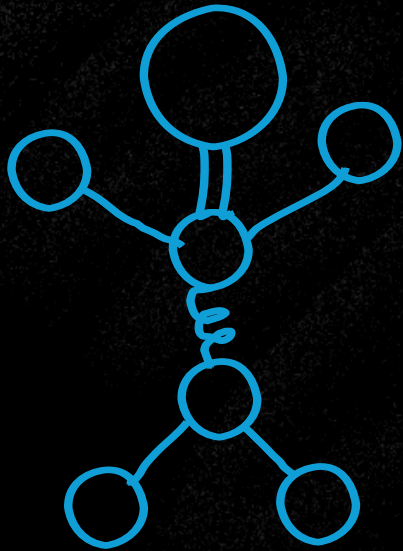Softmax

Scale

Matmul

Q    K    V

Transformer

# Fineprints of Attention – Bias

- positional and structural encodings
- biased attention
- lots and lots of long-range data

Matmul

Softmax

Scale

Matmul

E    Q    K    V

Transformer

# Fineprints of Attention – Data

- positional and structural encodings
- biased attention
- lots and lots of long-range data



Do I need to aggregate from distant neighbors?

Matmul

Softmax

Scale

Matmul

E    Q    K    V

Transformer

# Attention — Position and Structure

- WL Expressivity — why simple GNNs are not enough

- Positional encodings via eigenvectors

- Structural encodings via random walks

- Relative positions via distances and heat kernels

# WL Expressivity 🕶️

Weisfeiler-Lehman

# WL Expressivity 👓

Let's play a game!
Are these graphs the same?
Let's shrink again and count the neighbors
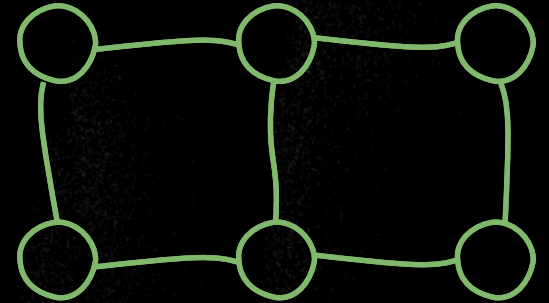
# WL Expressivity 🕶

# WL Expressivity 👓

They're the same!
Wait... That's not right.

There are things we cannot see
from inside because we do not
have position or direction!

I know what you're thinking.
Graphs have no direction

We'll circle back on that...
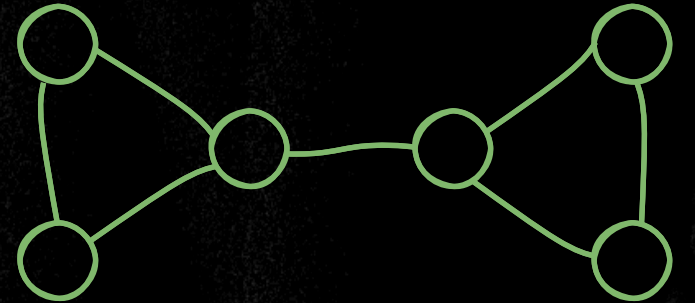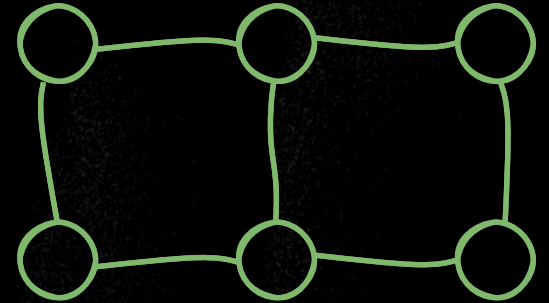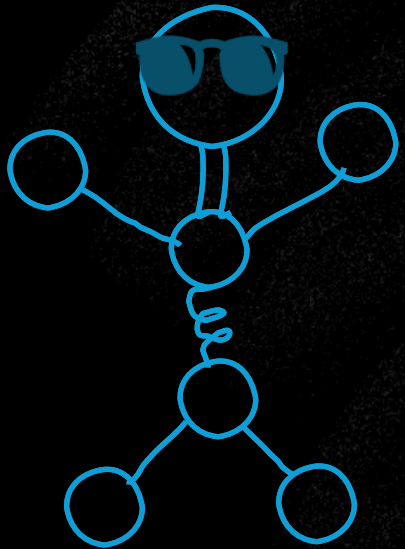
2  2
3  3
2  2
2  2
3  3
2  2
3  3
3  3
2  2

# Higher order features

Can we find some higher-order features?

Features that are permutation invariant but that can be computed?

Perhaps features inspired by higher order WL-tests by walking around the graph?

Let's look at random walks and motif detection
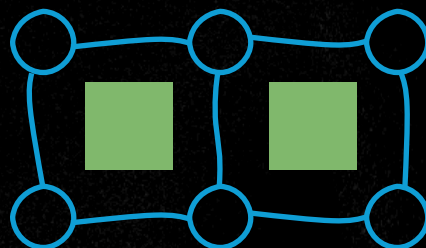
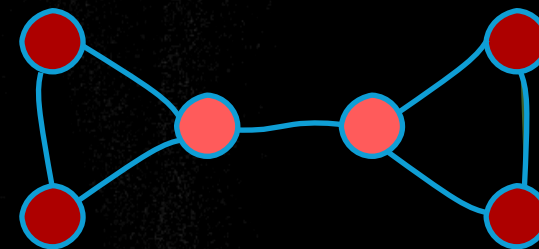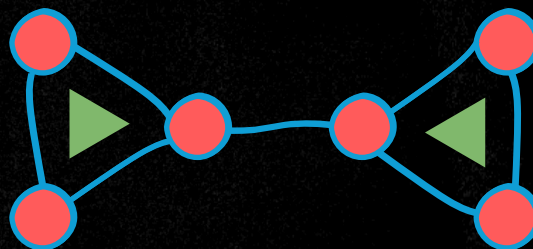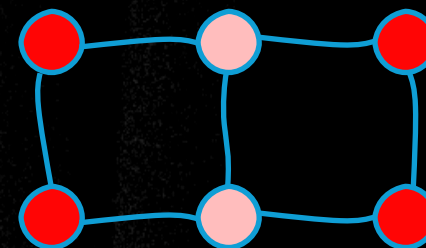# Structural encodings 🎡

These are nice local encodings!

But is there anything more global?

Random walk 3-step

Random walk 4-step

Now we can distinguish the graphs and nodes!
We can concatenate them to node features
We can bias the connectivity of the message passing
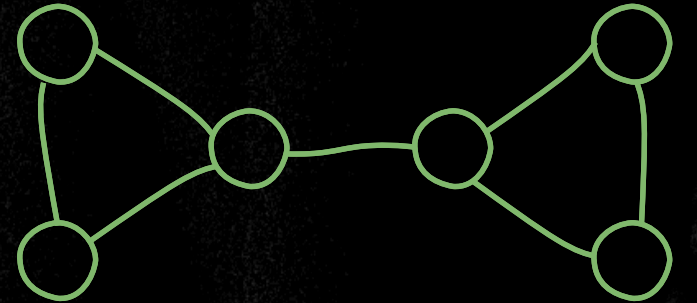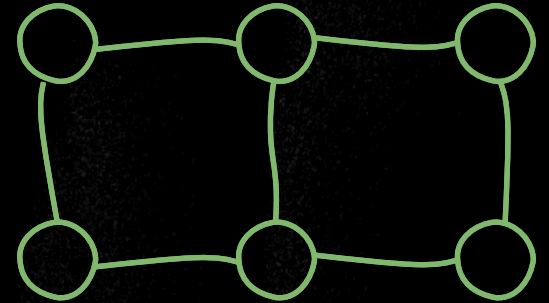We are again more expressive
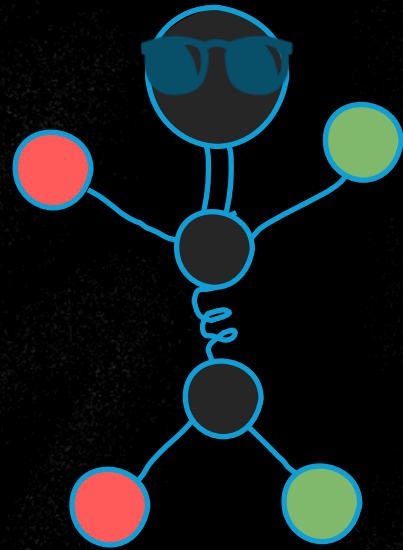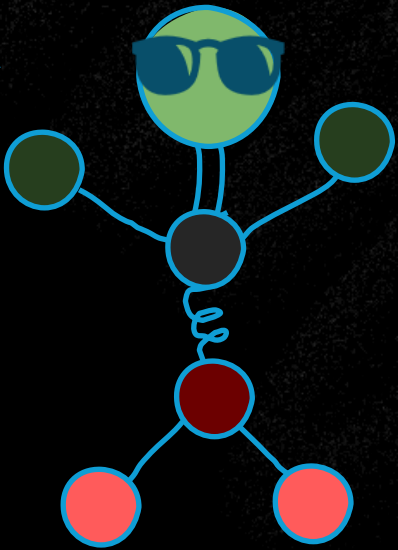
# Positional encodings 🗺

Low-frequency eigenvectors of the Laplacian
(lowest non-0 eigenvalue)

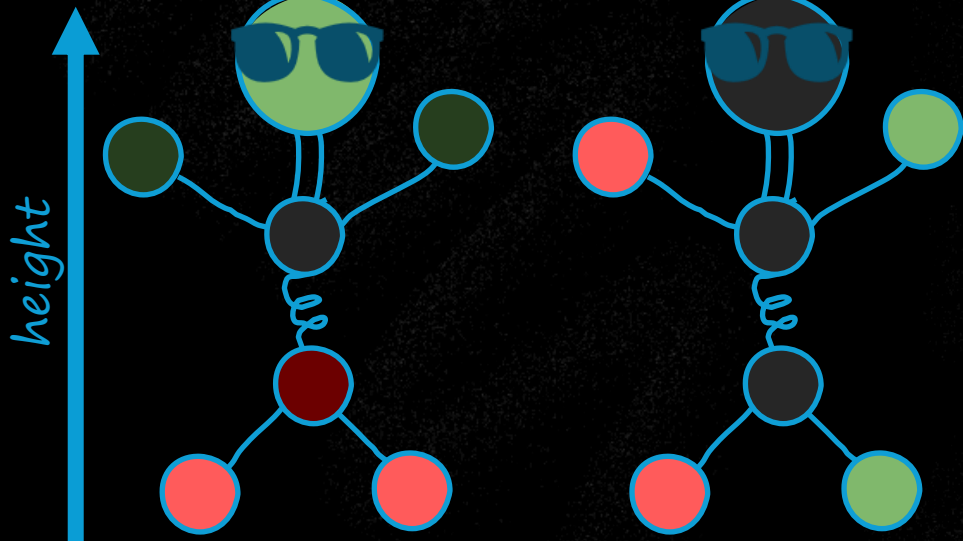$\phi_1$         $\phi_2$

Width

height

# Positional encodings 🗺

Low-frequency eigenvectors of the Laplacian
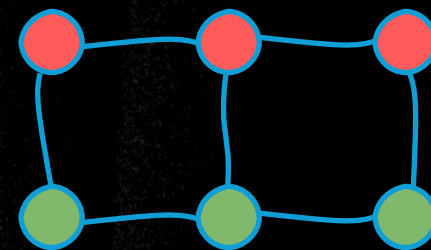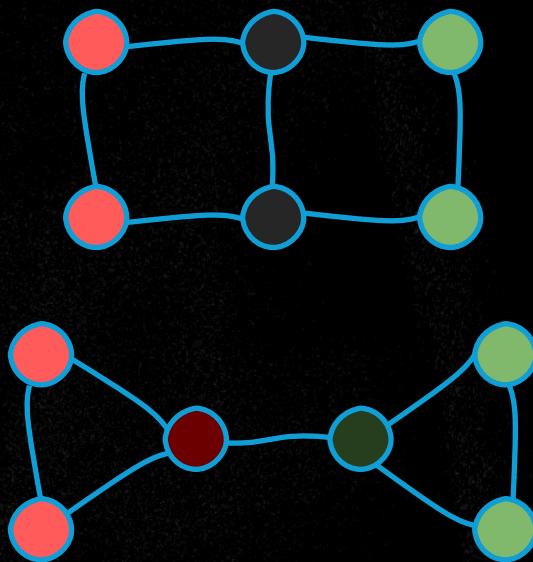(lowest non-0 eigenvalue)



$\phi_1$          $\phi_2$

Width

height

$\phi_1$

$\phi_2$

Longest length      Second longest length

Now we can distinguish the graphs and nodes!
We can concatenate them to node features
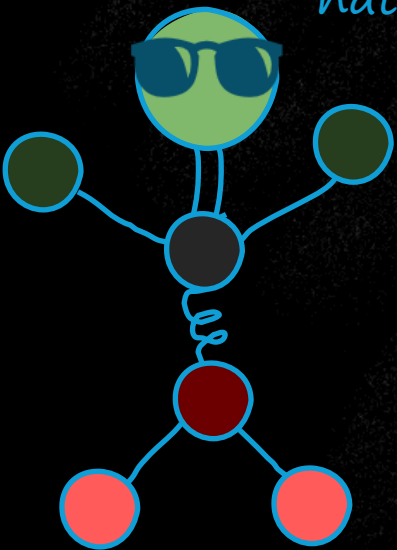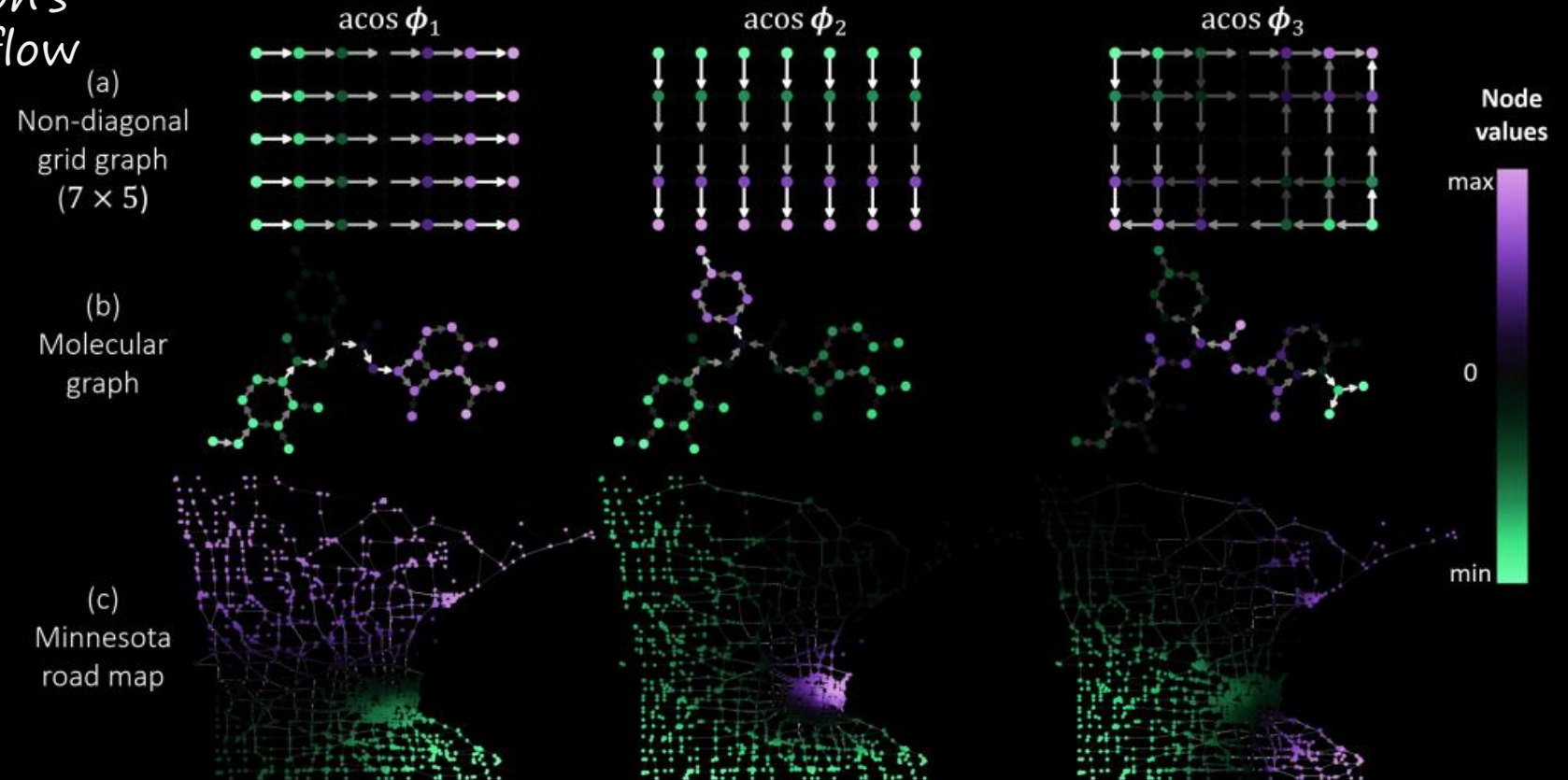We can bias the direction of the message passing
We are more expressive

# Low-frequency eigenvectors

The DGN work showed that they generalize CNNs when applied to grid graphs

They retrieve a graph's natural directional flow

Directional Graph Networks

Examples of eigenvector-based directions

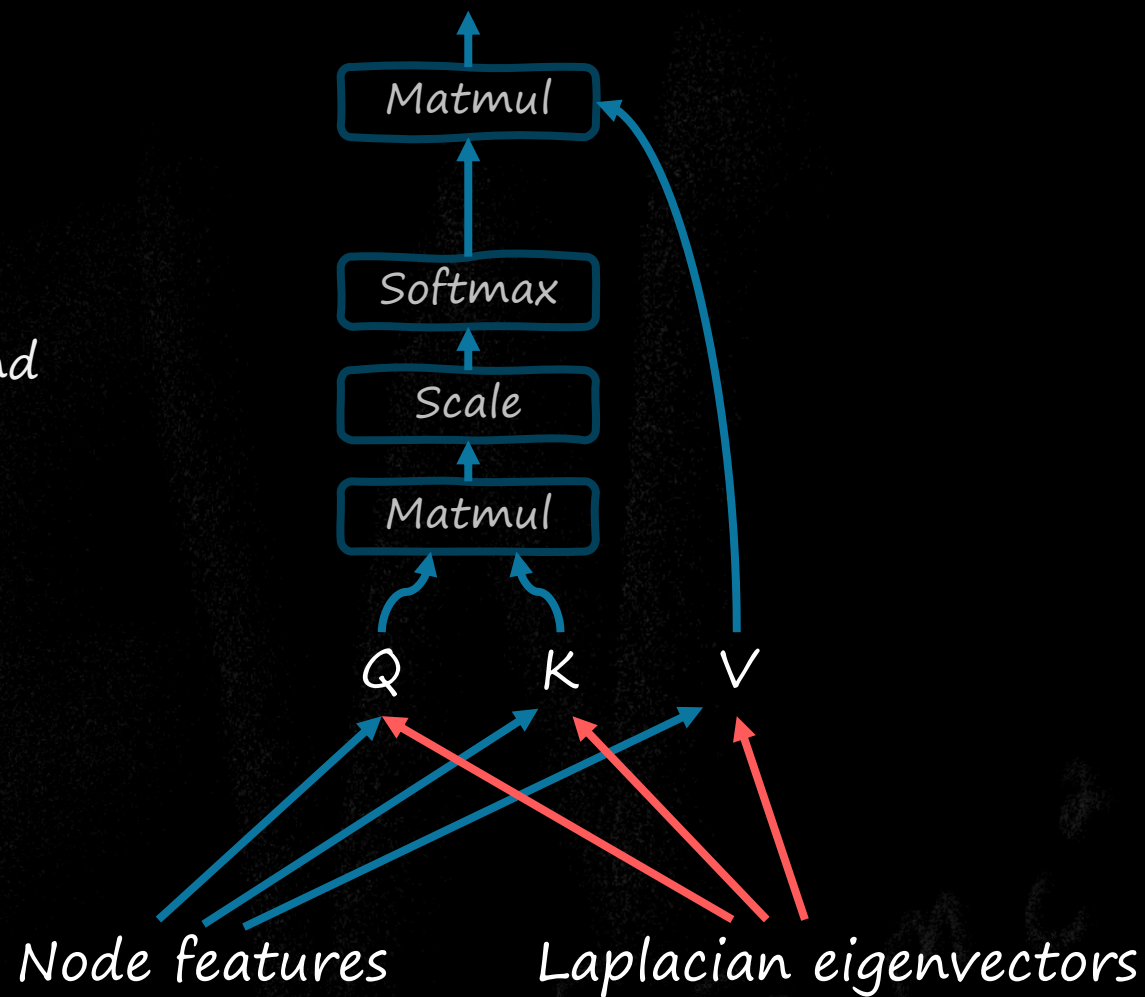$\text{acos } \phi_1$   $\text{acos } \phi_2$   $\text{acos } \phi_3$

(a) Non-diagonal grid graph $(7 \times 5)$

(b) Molecular graph

(c) Minnesota road map

Node values

max

0

min

# Basic graph Transformers

Basic graph Transformers have very poor results

- The connectivity is a strong inductive bias
- The eigenvectors are noisy and hard to understand for the network
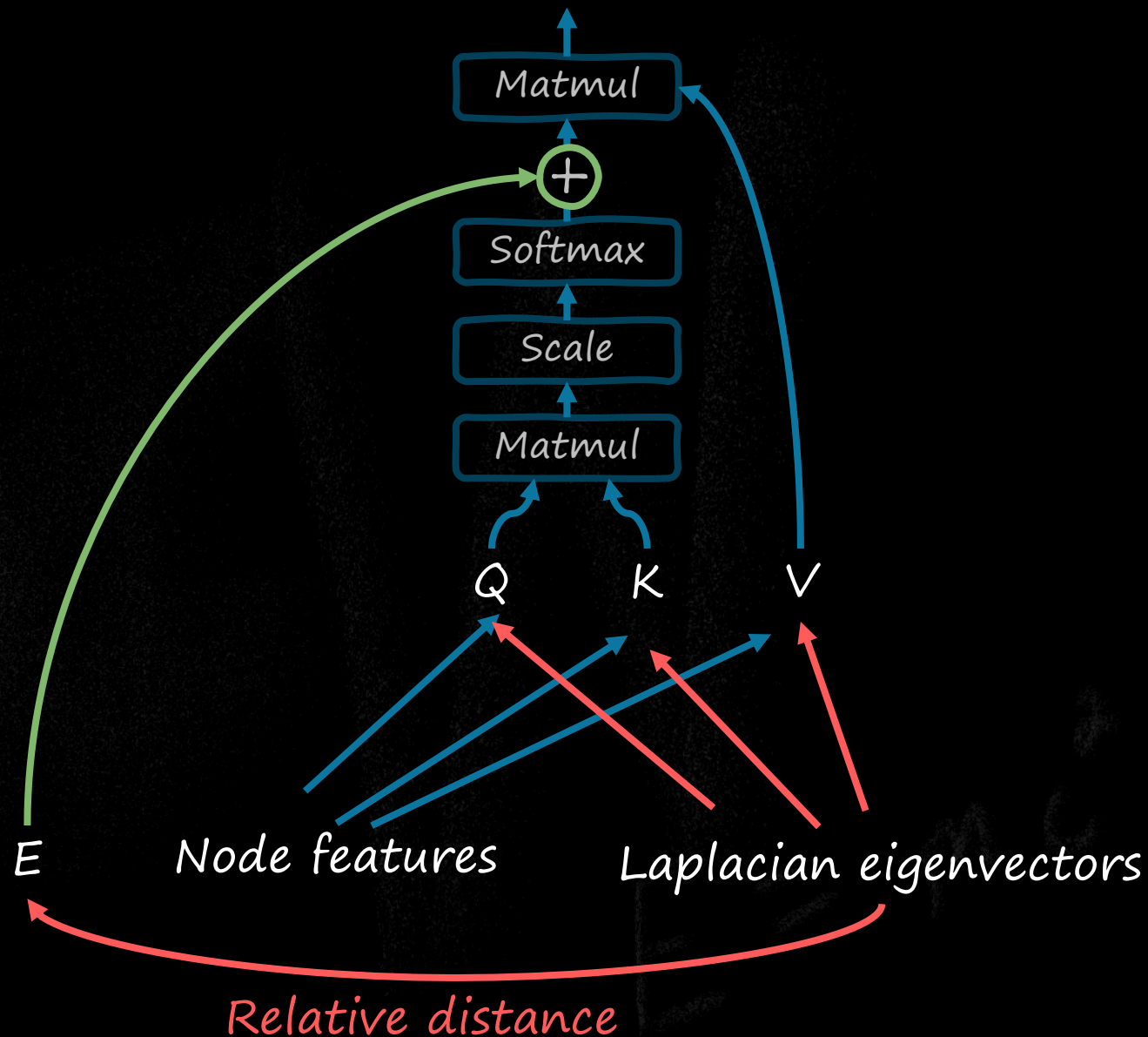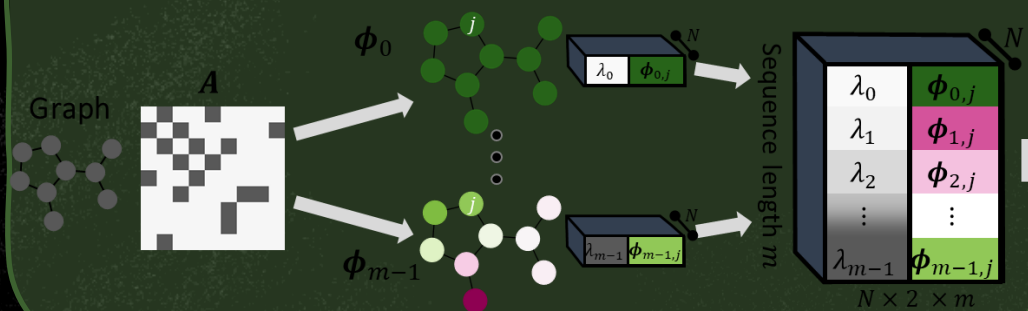- Edge features are missing

Matmul

Softmax

Scale

Matmul

Q        K        V

Node features                Laplacian eigenvectors

# Biased Full-Attention 🕸️



Thank you Bias!

Real edge?

Matmul

Softmax

Scale

Matmul

Q  K  V

E

Node features
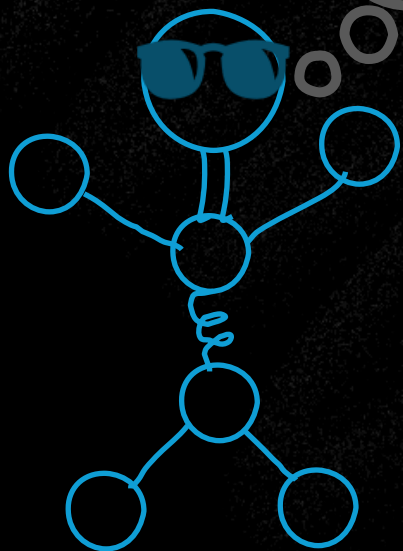
Laplacian eigenvectors

Relative distance

# Pre-training 🏋️

*Since we need
lots and lots of data,
Let's do pre-training.*

*How do we pre-train a
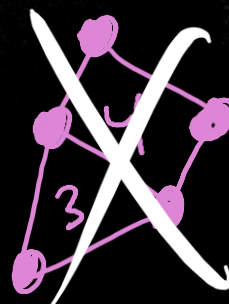molecular representation?*

## Biology
Cell assays
Transcriptomics

## Chemistry
Protein assays
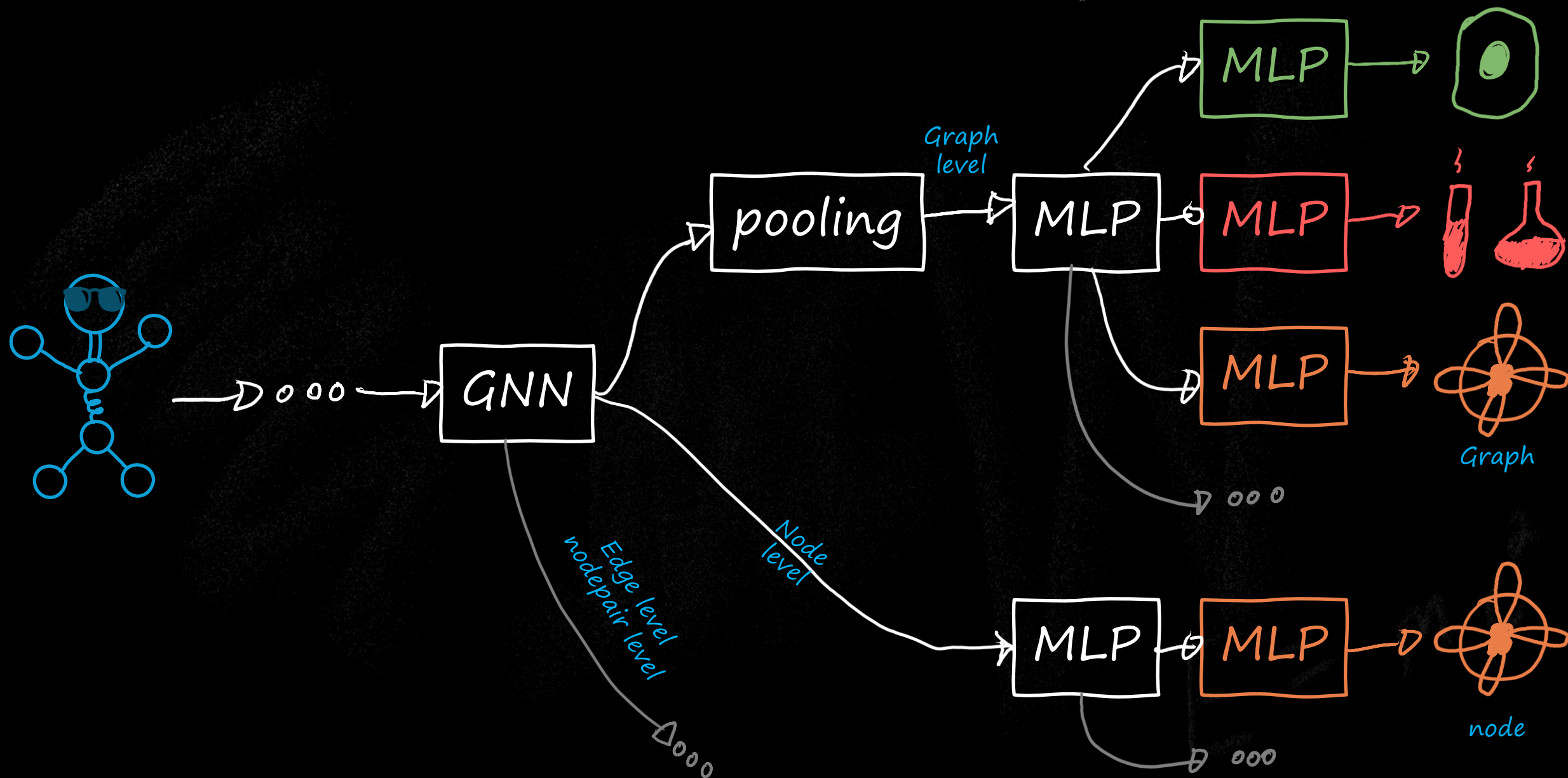Physicochemical (solubility, etc.)
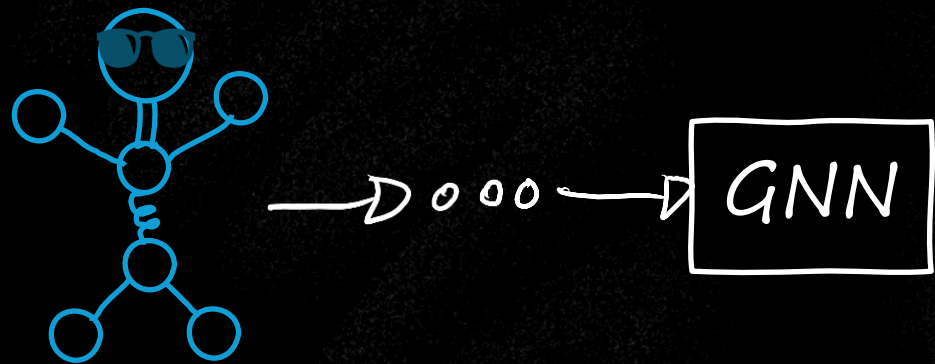
## Quantum mechanics
HOMO-LUMO gap
Partial charges

## Self-supervision
Finding missing atoms
Enumerating structures
SMILES reconstruction

# Multi-level multi-tasking 🗼

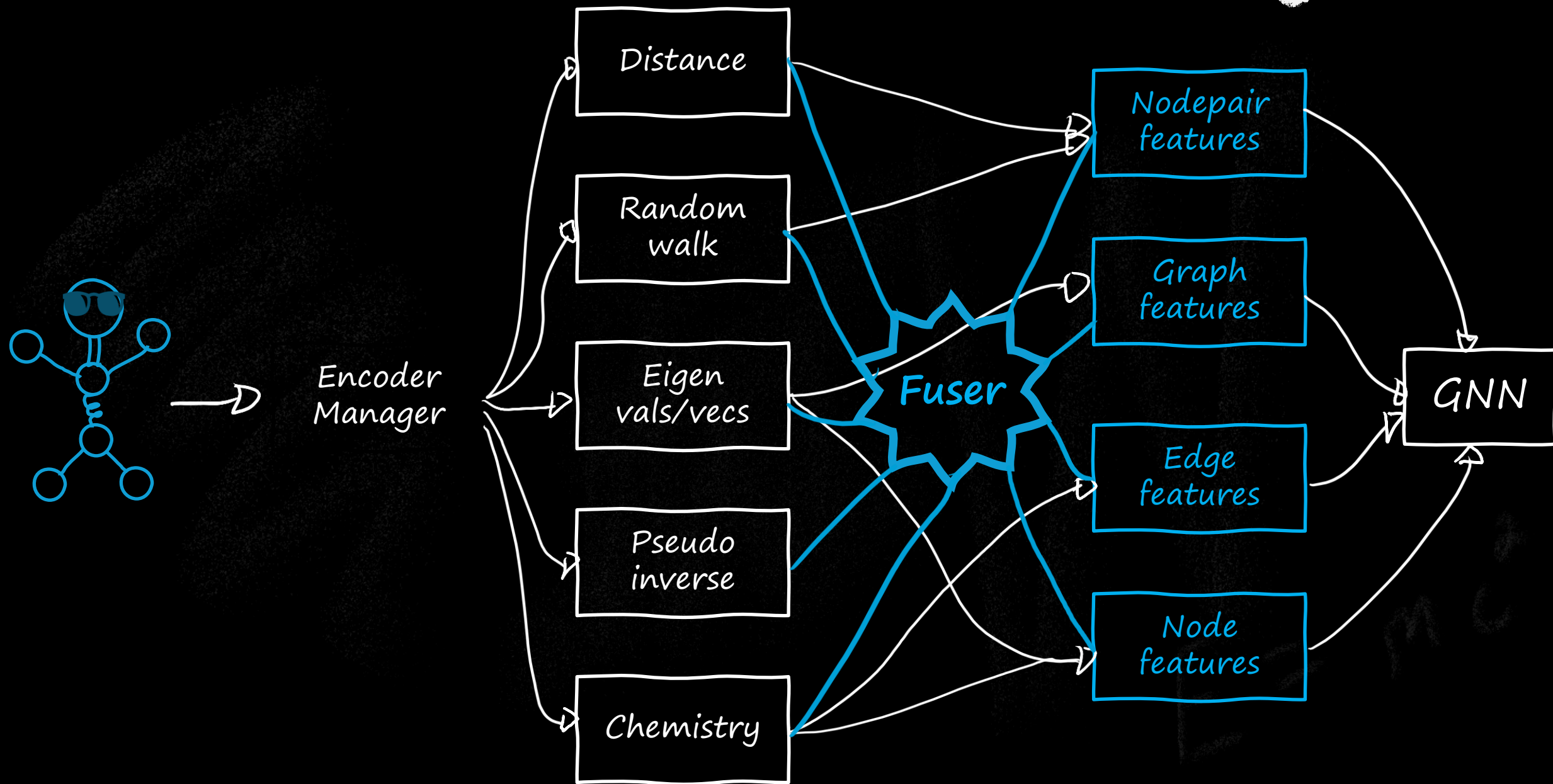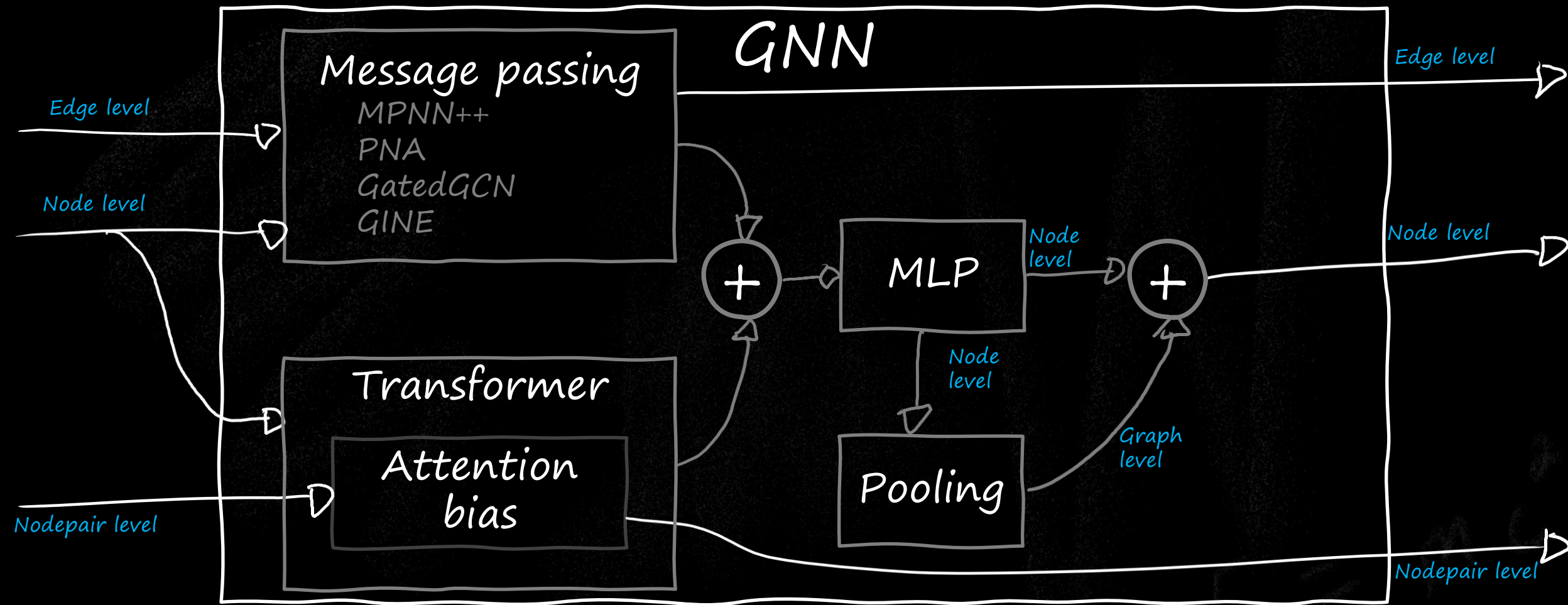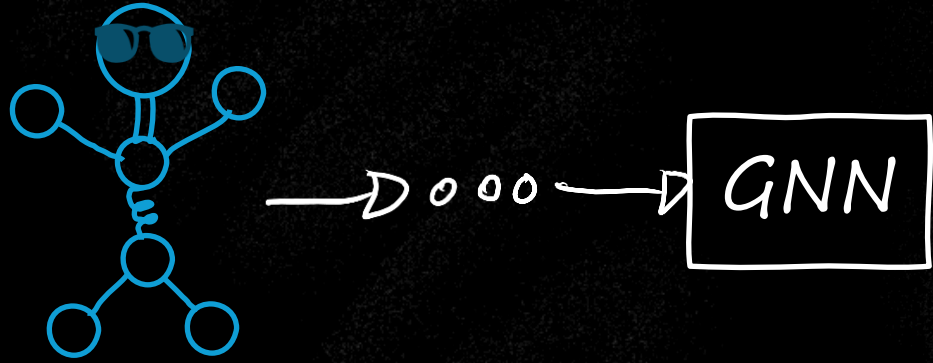# Multi-level multi-tasking 🗼

GNN

# Multi-Level positional encoding

# Multi-Level graph Transformer 🗼

# Finetuning 🏹

GNN

# Finetuning ⛄



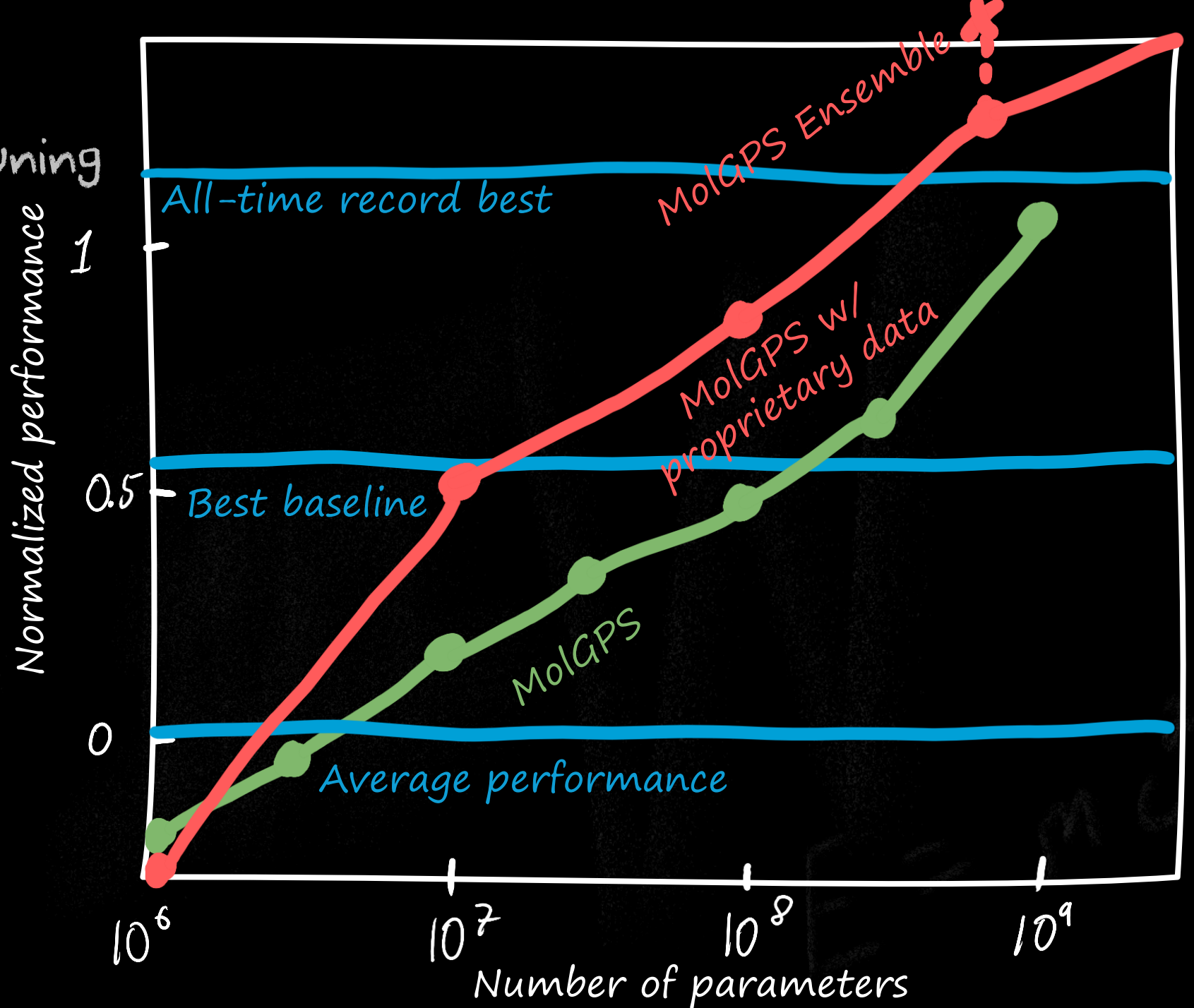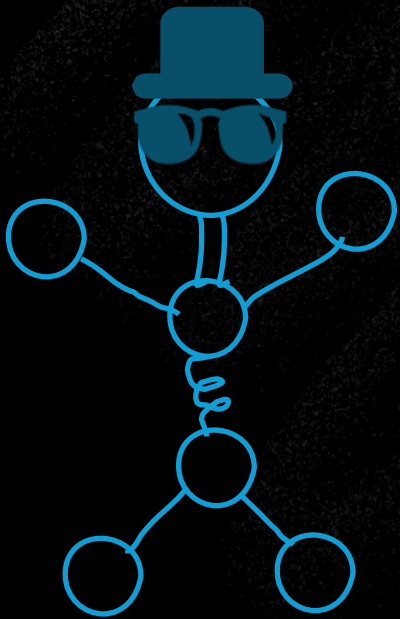New task

Graph level

pooling

MLP

GNN

Edge level
nodepair level

Node
level

MLP

MLP

MLP

MLP

MLP

MLP

Graph

node

# Limitation in no-context finetuning

- Even a model that perfectly understand physics and biology will overfit without context of the task

- Tasks can be encoded as
  - natural language
  - protein sequences
  - Cellular context

- Multimodality allows to encode context, and will make the GNNs much much more powerful

# Phenomics screening of molecules



Perturbation

1μM  5μM  0.1μM

Phenomics

How to build a model for this **equation** to understand how **molecules** impact **cells**?

# Let's try some contrastive Learning

Contrastive Loss



Prior methods have not succeeded, achieving only 8% recall

# The 3 Challenges
## and The Piano in New York

Is it visually ~~inactive?~~

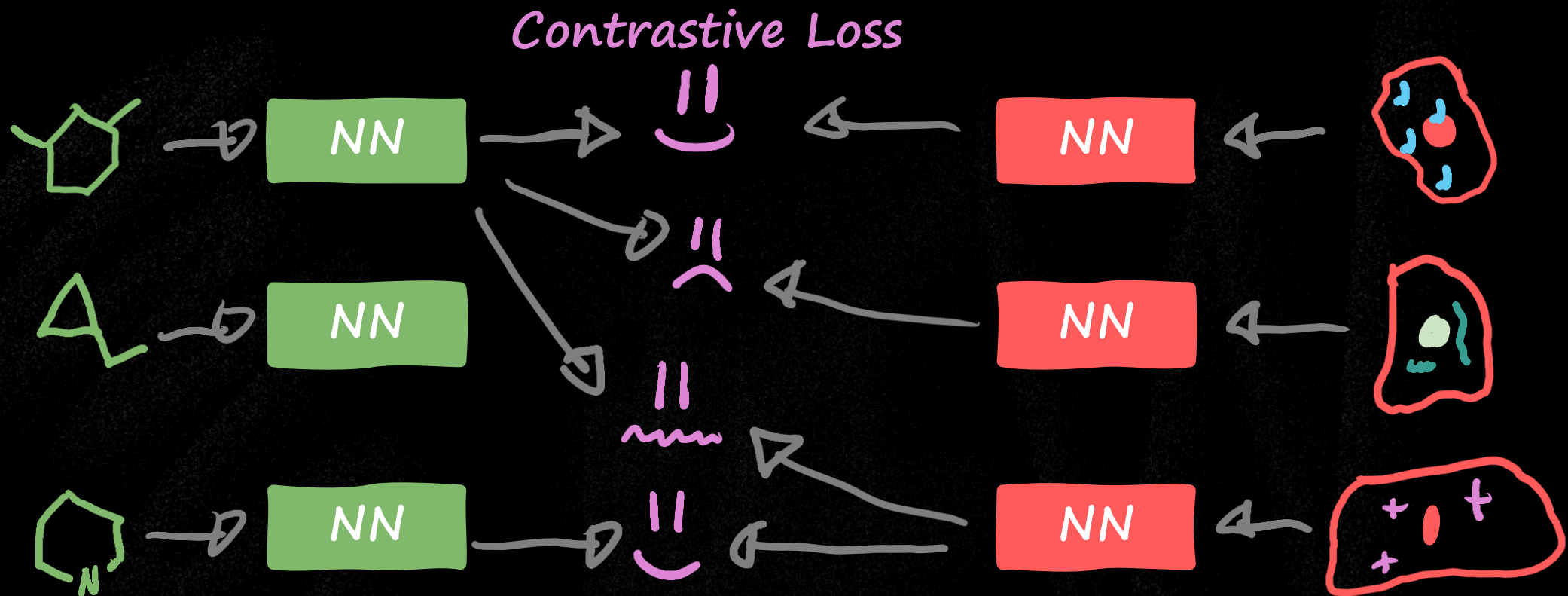Background noise

Which ~~molecule~~ is playing?

music

Increase the ~~concentration~~

volume

7pm

4pm

# The 3 Challenges

**Natural variations:** Batch effects are the largest source of variation in phenomics images.

Can we ignore it?

**90% inactives:** Most molecules have no visible effect.

How to handle this source of random noise?

**Concentration**
Too low: Nothing happens
Too high: Everyone dies

How to model this non-linear relationship?

# Better contrastive Learning with S2L Loss
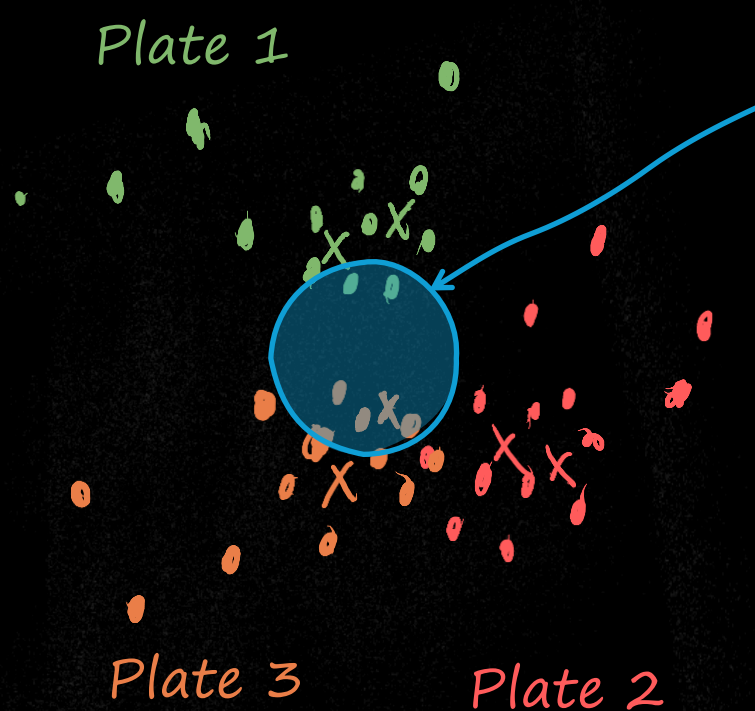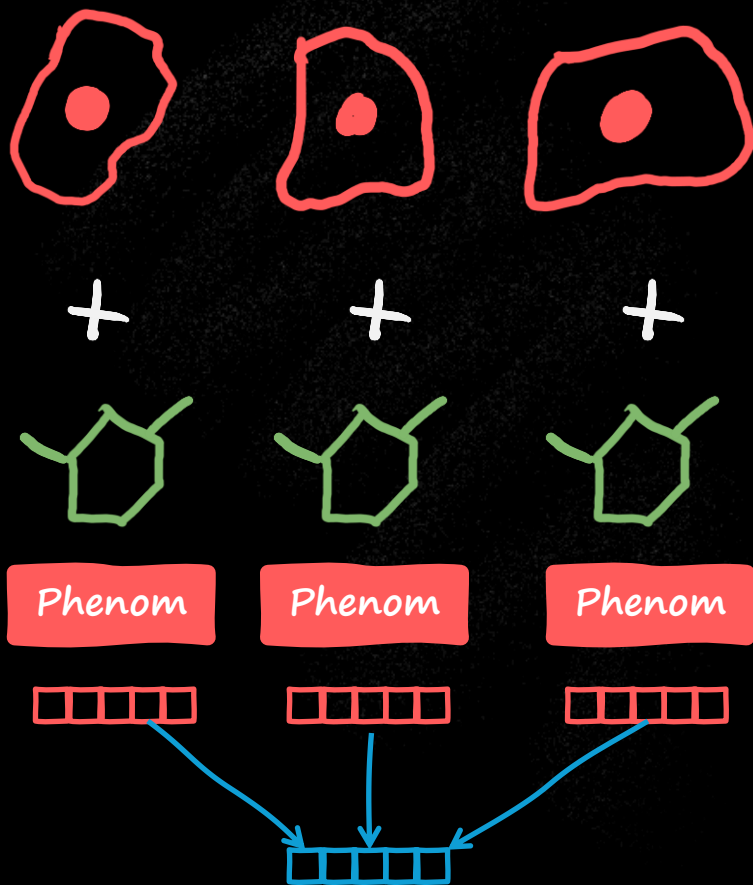
Molecular perturbation

MolGPS — Pretrained

Explicit Concentration → MLP

S2L Contrastive Loss ‿ ⁀

Implicit Concentration → MLP

Phenom — Pretrained

Experimental image

Re-weight samples based on Phenomics embedding
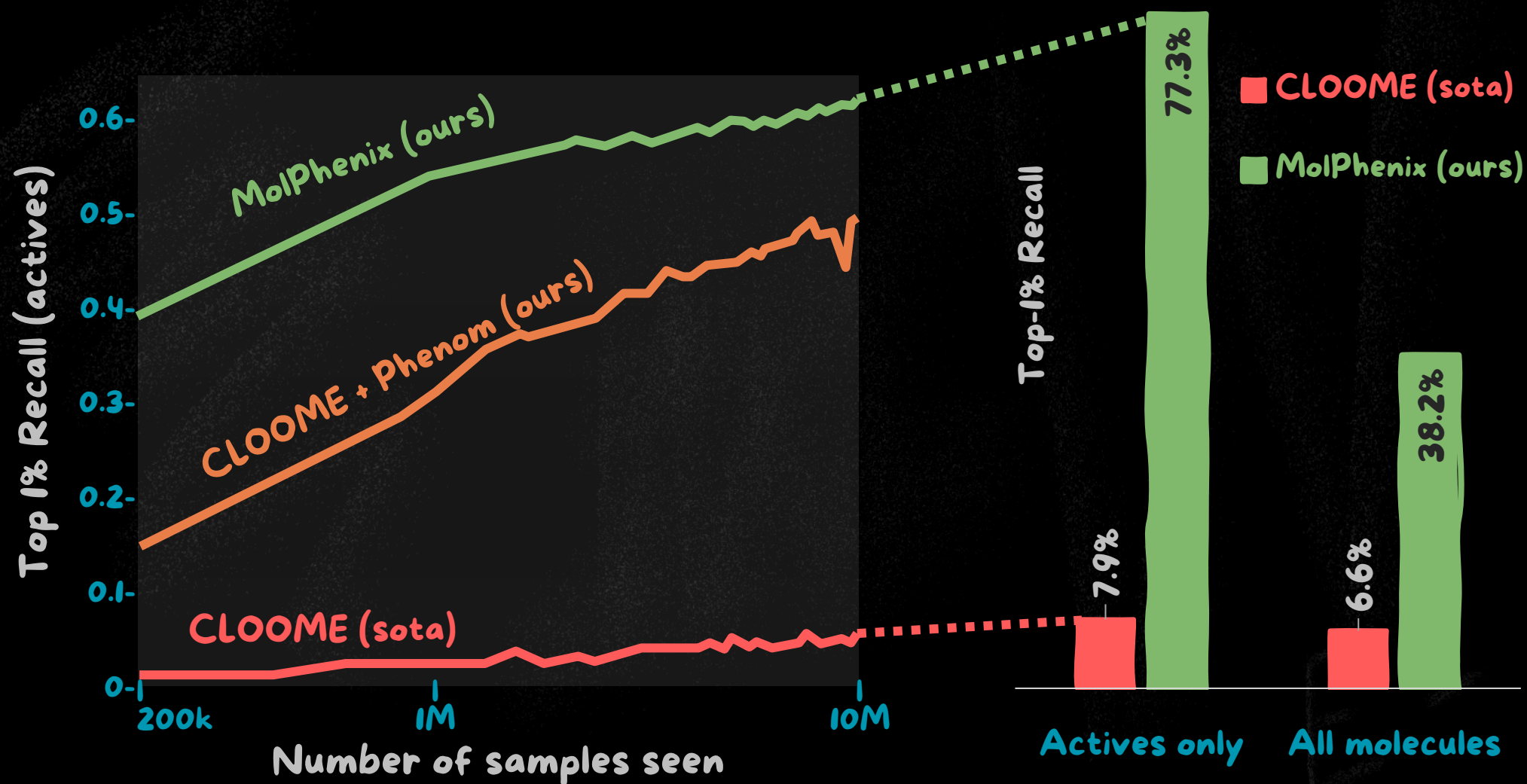
$$\mathcal{L}_{S2L} = -\frac{1}{N}\sum_{i-1}^{N}\sum_{j-1}^{N}\log\left[\frac{w_{i,j}^{x}}{1+\exp(-\alpha\langle \boldsymbol{z_X}, \boldsymbol{x_m}\rangle + b)} + \frac{1 - w_{i,j}^{x}}{1+\exp(\alpha\langle \boldsymbol{z_X}, \boldsymbol{x_m}\rangle + b)}\right]$$

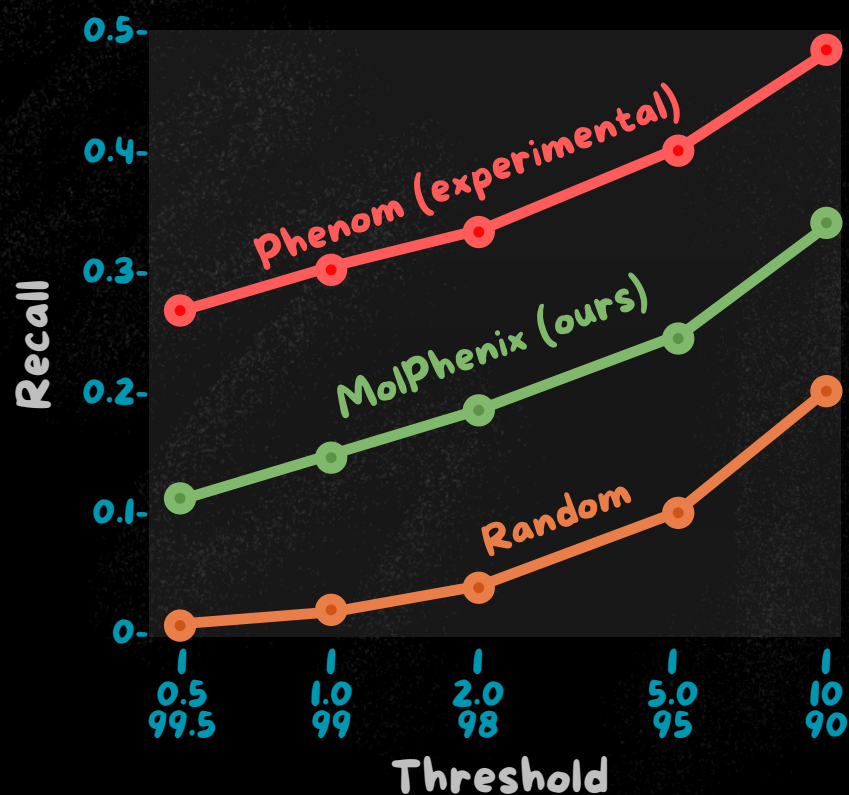Use sigmoid to reduce the effect of false negatives

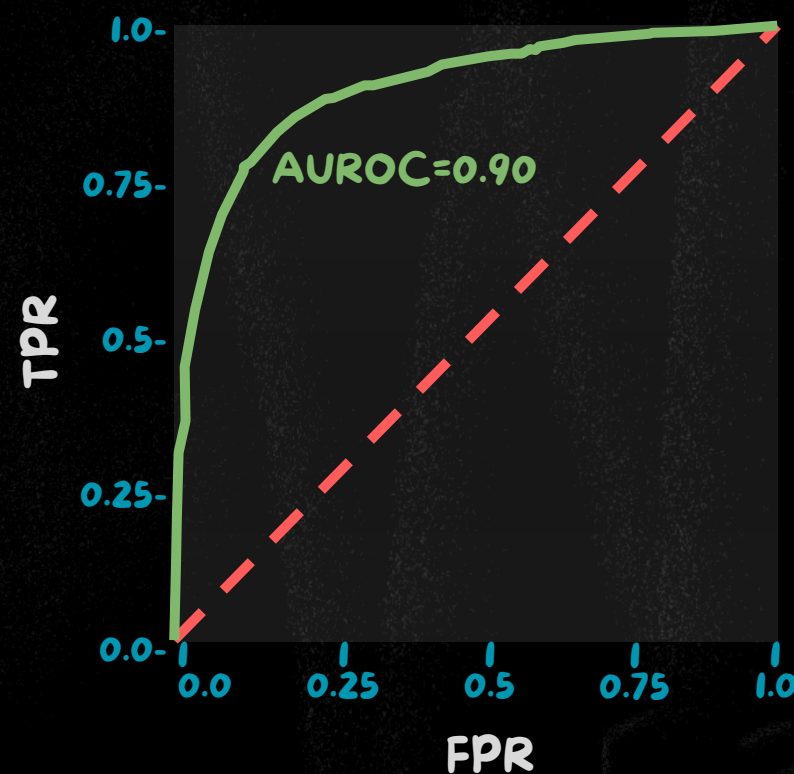# 10x recall compared to previous SOTA
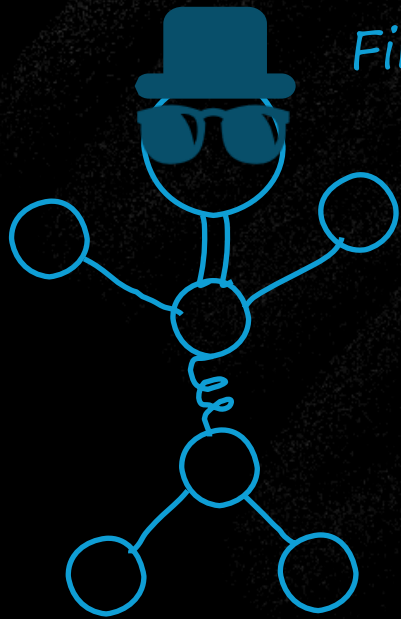
# Downstream applications

# MolPhenix

- MolPhenix opens a completely new direction for ML in drug discovery with 10x improvement

- We can model how molecules impact cells, not just do some predictive assays →

- A first step towards Virtual Cells, to industrialize drug discovery in the age of AI

# Thank you Graphy!! And Dom

Finally, Dom will stop talking!

But if you're not tired of him, you can follow him on ~~Twitter~~ X
@Dom_Beaini

Thanks to a thousand co-authors!