

# PREMISE SELECTION FOR HIGHER-ORDER ATP PROBLEMS

Keneni W. Tesema, Jan Jakubuv, Martin Suda

Czech Technical University in Prague

Deeper4AI Project

Gothenburg CHAIR workshop, 28<sup>th</sup> May 2026

# INTRO: PREMISE SELECTION

- Premise selection is the task of identifying, from a library of potentially thousands of axioms, a small subset necessary to prove a given conjecture.



- ATP without premise selection



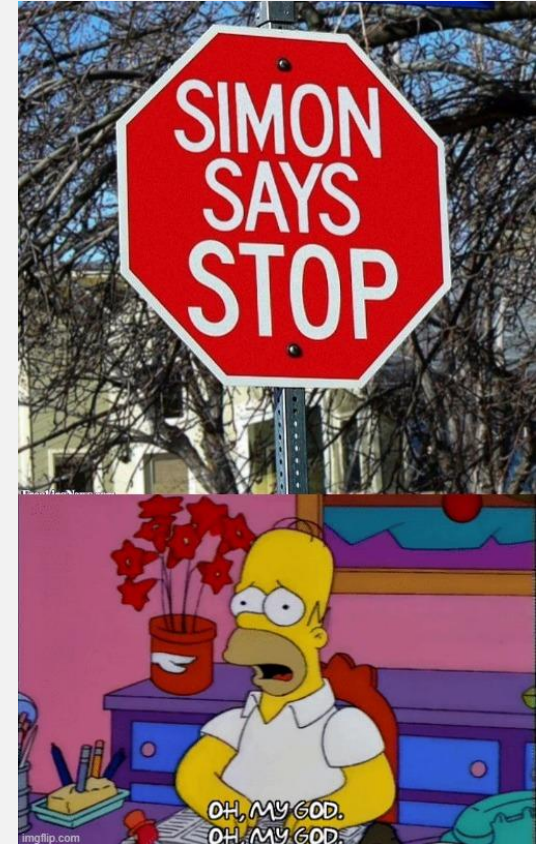
- ATP with premise selection



- Magic

# INTRO: HIGER-ORDER LOGIC (HOL) PROBLEMS

- **First-order:** "Be confident."
- **Higher-order:** "Be confident in your ability to become more confident."
- **First-order:** "Translate 'hello' to 'bonjour'."
- **Higher-order:** "Translate any instruction from English to French, including the word 'translate' itself."
- **First-order:** "Do what Simon says"
- **Higher-order:** "Do what Simon says when Simon says "Simon says....""



# PREMISE SELECTION IN HOL

FOL	HOL
<b>Variables</b> range over <b>individuals</b> (people, numbers, etc.)	<b>Variables</b> can range over <b>FUNCTIONS</b> (properties etc): $\forall P. P(\text{agatha}) \rightarrow P(\text{butler})$
<b>Predicates</b> apply to individuals: hates (agatha, X)	<b>Lambda abstraction:</b> $\lambda x. x + x$ (functions as data)
<b>Functions</b> return individuals: father_of (agatha)	<b>Partial application:</b> hates(agatha) returns a function $\lambda Y. \text{hates}(\text{agatha}, Y)$
<b>Quantifiers</b> only over individuals: $\forall X. \text{hates}(\text{agatha}, X)$	<b>Higher-order unification:</b> Matching $P(\text{agatha})$ against $\text{hates}(\text{agatha}, \text{butler})$ requires finding $P = \lambda X. \text{hates}(X, \text{butler})$
Simple Abstract Syntax Trees(AST) Eg: Ax: hates(agatha, X) <div style="text-align: center; margin-top: 20px;"> <pre>               hates              /   \             /     \           agatha   X           </pre> </div>	Richer AST Eg: $^{\lambda X:\$i} : (@(@(\text{hates}, \text{agatha}), X))$ <div style="text-align: center; margin-top: 20px;"> <pre>               λ                               X                               @              / \             /   \            @     X           / \          /   \         /     \        /       \       /         \      /           \     /             \    /               \   /                 \  /                   \ /                     \ agatha             hates           </pre> </div>

# CHALLENGES IN PREMISE SELECTION IN HOL

FOL

HOL

## Challenge 1: Variable Arity and Partial Application

hates(agatha, butler)

hates(agatha), (hates(agatha))(butler), @(@(hates, agatha), butler)

## Challenge 2: Lambda Abstractions (Anonymous functions)

None

$\lambda x. x + x$

## Challenge 3: Rich Type system

Simple types (i, o)

Function types, Polymorphic types, Type variables and instantiation

## Challenge 4: Many Node Types

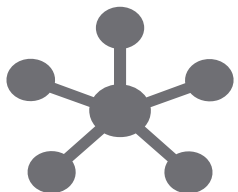
Symbol, variable, operator

Application (@), Lambda( $\wedge$ ), Type abstraction, Type application, implicit vs explicit quantifications

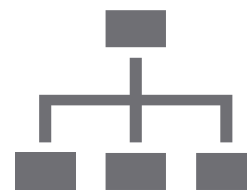
# RELATED WORKS

Technique	Approach	FOL	HOL	Type	Key Limitation
MePo (Meng & Paulson, 2009)	Symbol overlap heuristics	✓	✗	Heuristic	Shallow, no structure
MaSh (Kühlwein et al., 2013)	Naive Bayes + kNN ensemble	✓	✗	Heuristic	Hand-crafted features only
SInE (Hoder & Voronkov, 2011)	Least common symbol trigger (used in Vampire/E-prover)	✓	✗	Heuristic	Approximate, no types
Wang et al. (2017)	Graph embeddings with shared symbol nodes	✓	✓	GNN	First graph-based for HOL
Olšák et al. (2019)	Graph embeddings + property-invariant hypergraph	✓	✗	GNN	Invariant to symbol renaming; focuses on unseen Skolem symbols
ARCG-NN (Liu et al., 2021)	Cross-graph attention (conjecture ↔ axioms)	✓	✓	GNN	Better relevance encoding
Bansal et al. (2019)	GNN + sequence models (HOL Light)	✓	✓	Hybrid	HOL4/HOL Light specific
Magnushammer (Mikuła et al., 2024)	Contrastive transformer + InfoNCE (SELECT+RERANK)	✓	✓	Sequence	No explicit type structure, (types as atoms, not composed)
Paliwal et al. (2020)	Graph w/ types + binding	✓	✓	GNN	Explicit HOL structure
QUILL (Kogkalidis et al., 2024)	Structured attention + de Bruijn indices	✓	✓	GNN	Dependent types (Agda), not THF

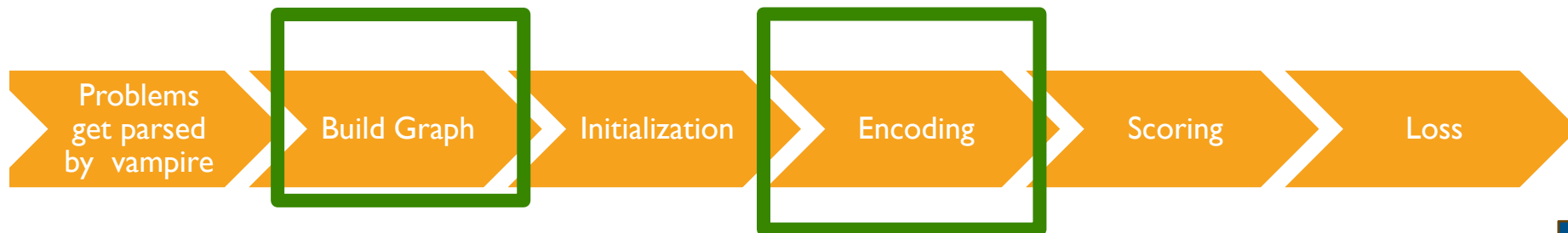
# METHODOLOGY



**Graph-based networks**



**Sequence-based networks**



Problem:  
conj,  
ax 1 to ax 2000

Problem:  
conj,  
ax 1 :4  
Ax 2: 7  
Ax3: -2  
Ax20 : 9

# EXAMPLE: KNIGHTS AND KNAVES

A very special island is inhabited only by knights and knaves. Knights always tell the truth, and knaves always lie. You meet two inhabitants: Zoey and Mel. Zoey tells you that Mel is a knave. Mel says, 'Neither Zoey nor I are knaves'. Who is a knight and who is a knave?

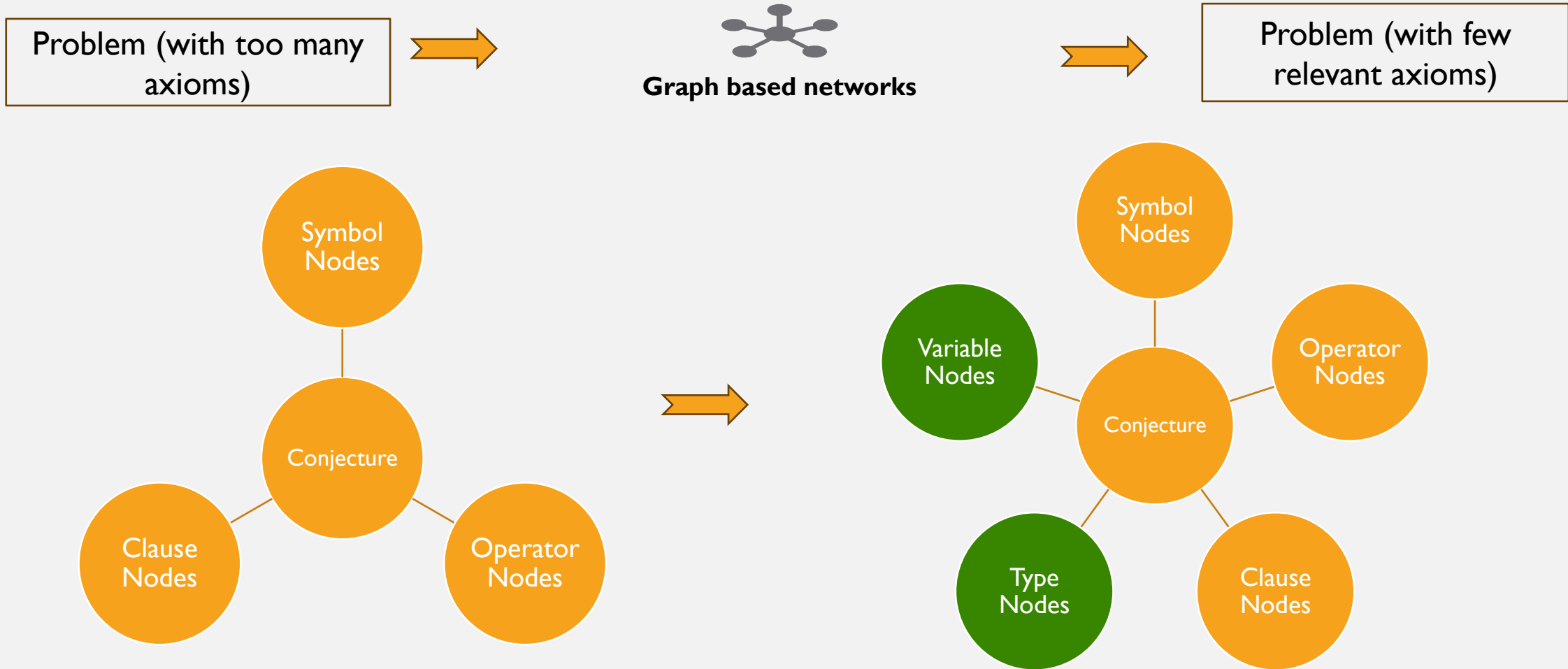
## Conjecture:

$\exists Y Z. ((Y = \text{knight}) \leftrightarrow (Y = \text{knight})) \wedge ((Z = \text{knight}) \leftrightarrow (Z = \text{knight})) \wedge \text{is\_a}(\text{mel}, Y) \wedge \text{is\_a}(\text{zoey}, Z)$

- **ax1:**  $\forall X. \text{is\_a}(X, \text{islander}) \rightarrow (\text{is\_a}(X, \text{knight}) \vee \text{is\_a}(X, \text{knave}))$   
(Islanders are only knights or knaves)
- **ax2:**  $\forall X. \text{is\_a}(X, \text{knight}) \rightarrow \forall A. (\text{says}(X, A) \rightarrow A)$   
(Knights tell the truth)
- **ax3:**  $\forall X. \text{is\_a}(X, \text{knave}) \rightarrow \forall A. (\text{says}(X, A) \rightarrow \neg A)$   
(Knaves lie)
- **ax4:**  $\text{is\_a}(\text{zoey}, \text{islander}) \wedge \text{is\_a}(\text{mel}, \text{islander})$   
(Zoey and Mel are islanders)
- **ax5:**  $\text{says}(\text{zoey}, \text{is\_a}(\text{mel}, \text{knave}))$   
(Zoey says Mel is a knave)
- **ax6:**  $\text{says}(\text{mel}, \neg(\text{is\_a}(\text{zoey}, \text{knave}) \vee \text{is\_a}(\text{mel}, \text{knave})))$   
(Mel says "Neither Zoey nor I are knaves")

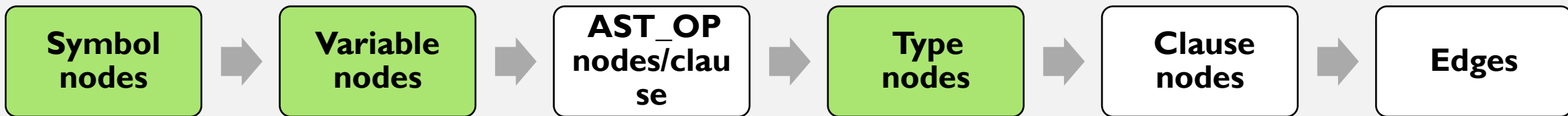


# METHODOLOGY: GNN\_VI

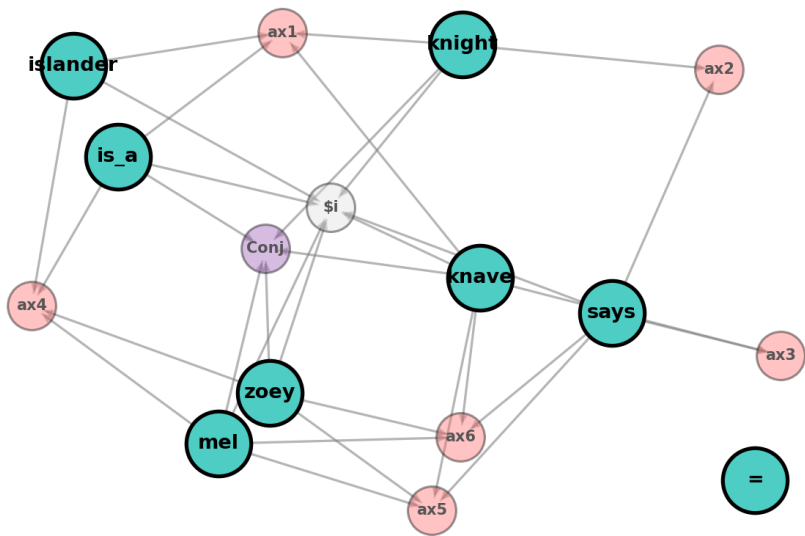


Olšák, Miroslav, Cezary Kaliszyk, and Josef Urban. "Property invariant embedding for automated reasoning." *arXiv preprint arXiv:1911.12073* (2019).

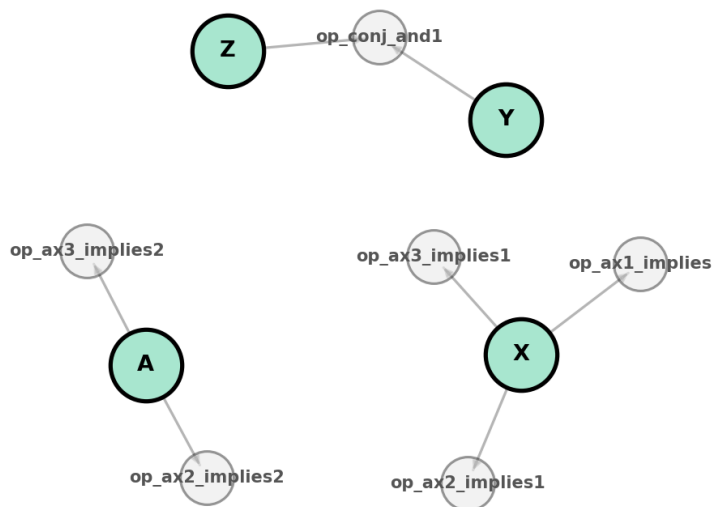
# GNN\_v1



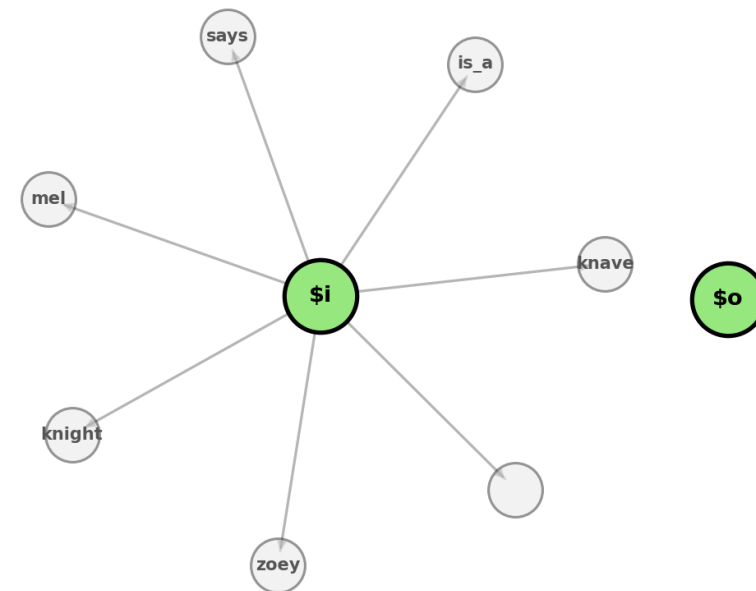
## Symbol Nodes



## Variable Nodes



## Type Nodes



types, sort\_terms, symbols, formula\_nodes, formula\_roles

# GNN\_v1

Symbol nodes



Variable nodes



AST\_OP nodes/clause



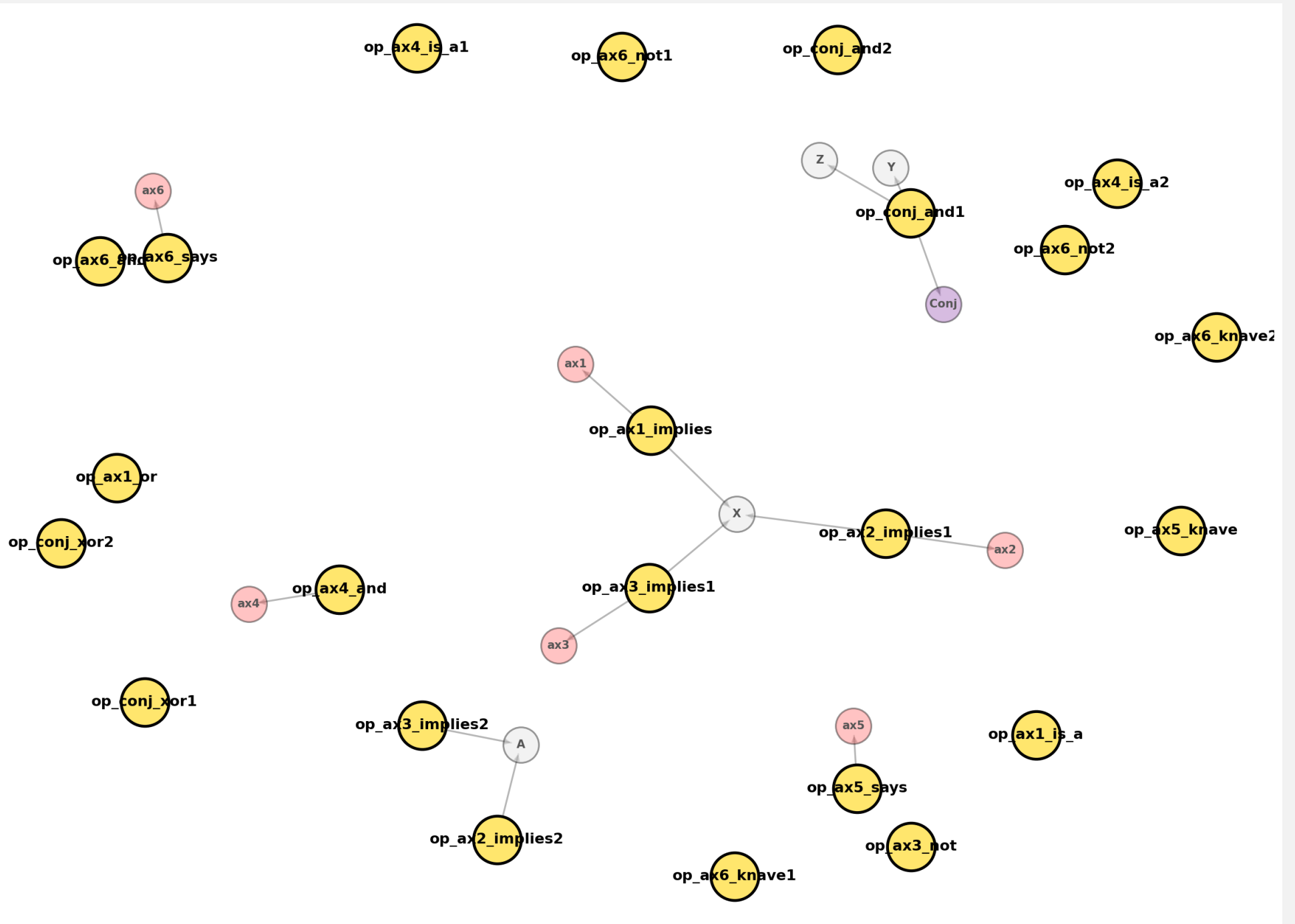
Type nodes



Clause nodes



Edges



# GNN\_v1

Symbol nodes



Variable nodes



AST\_OP nodes/clause



Type nodes

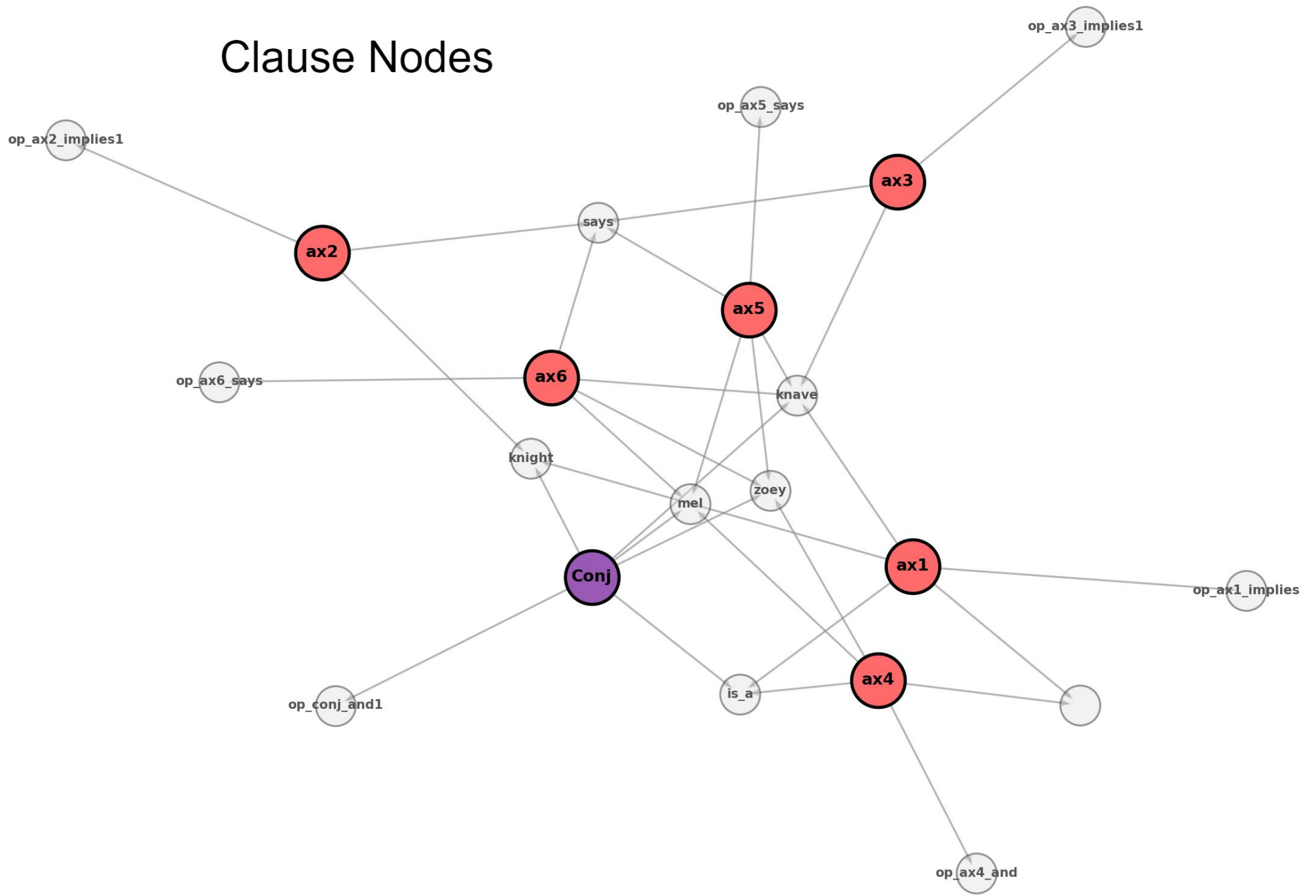


Clause nodes

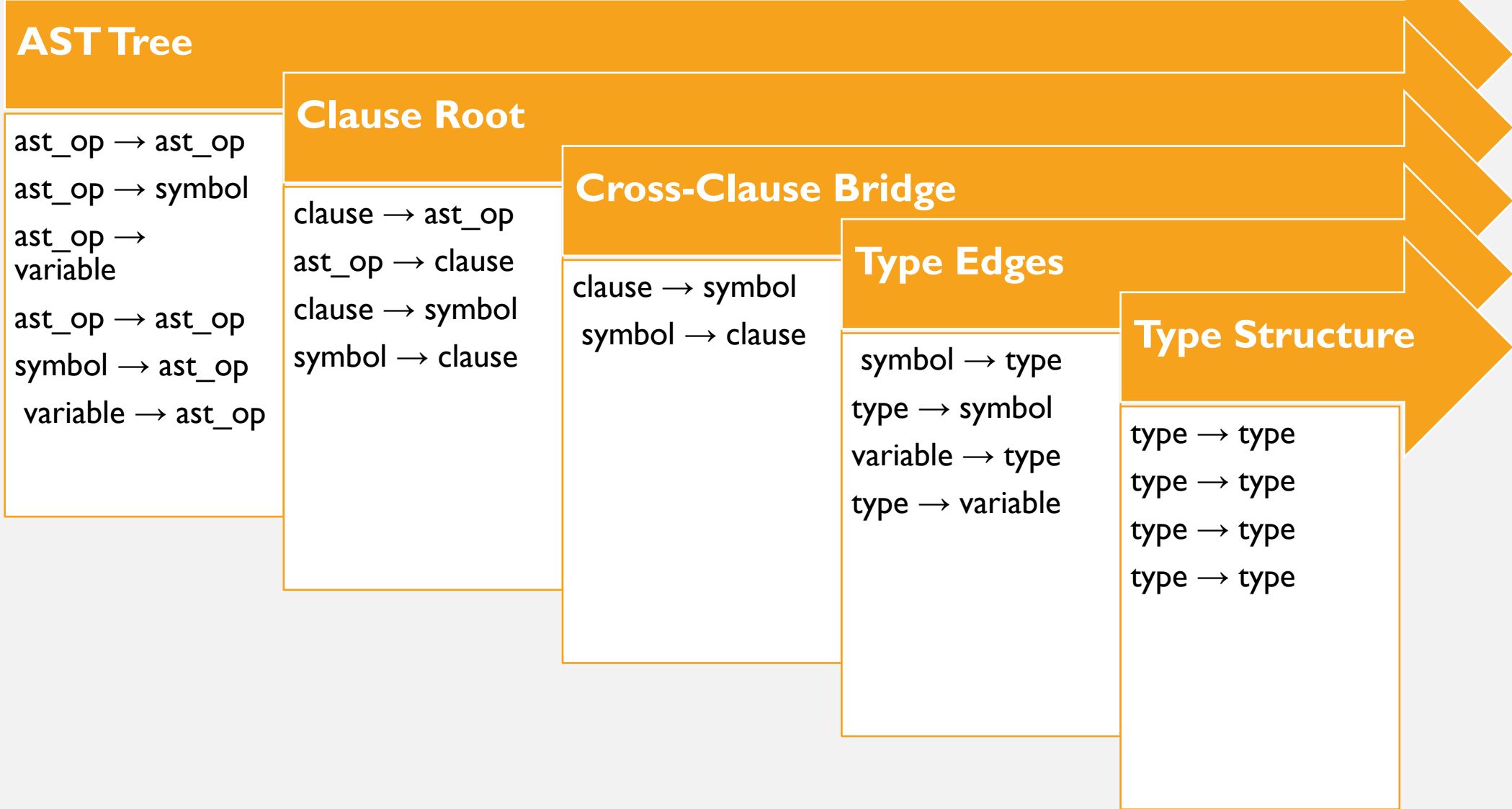
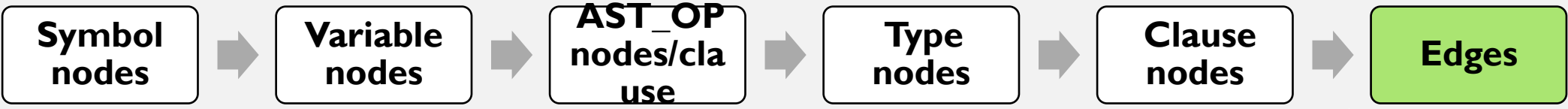


Edges

## Clause Nodes



# GNN\_v1



# GNN\_v1

Symbol nodes



Variable nodes



AST\_OP nodes/clause



Type nodes



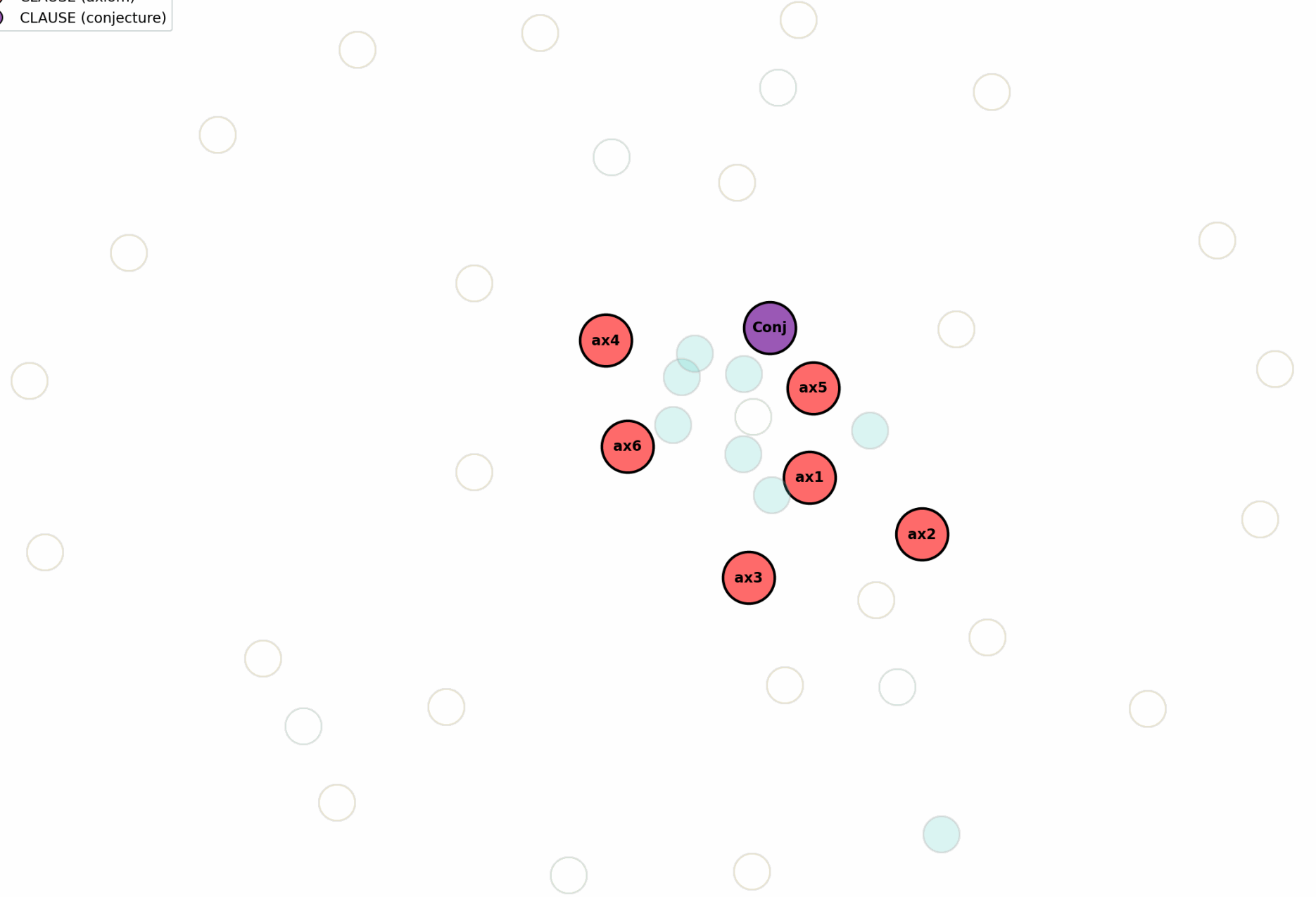
Clause nodes



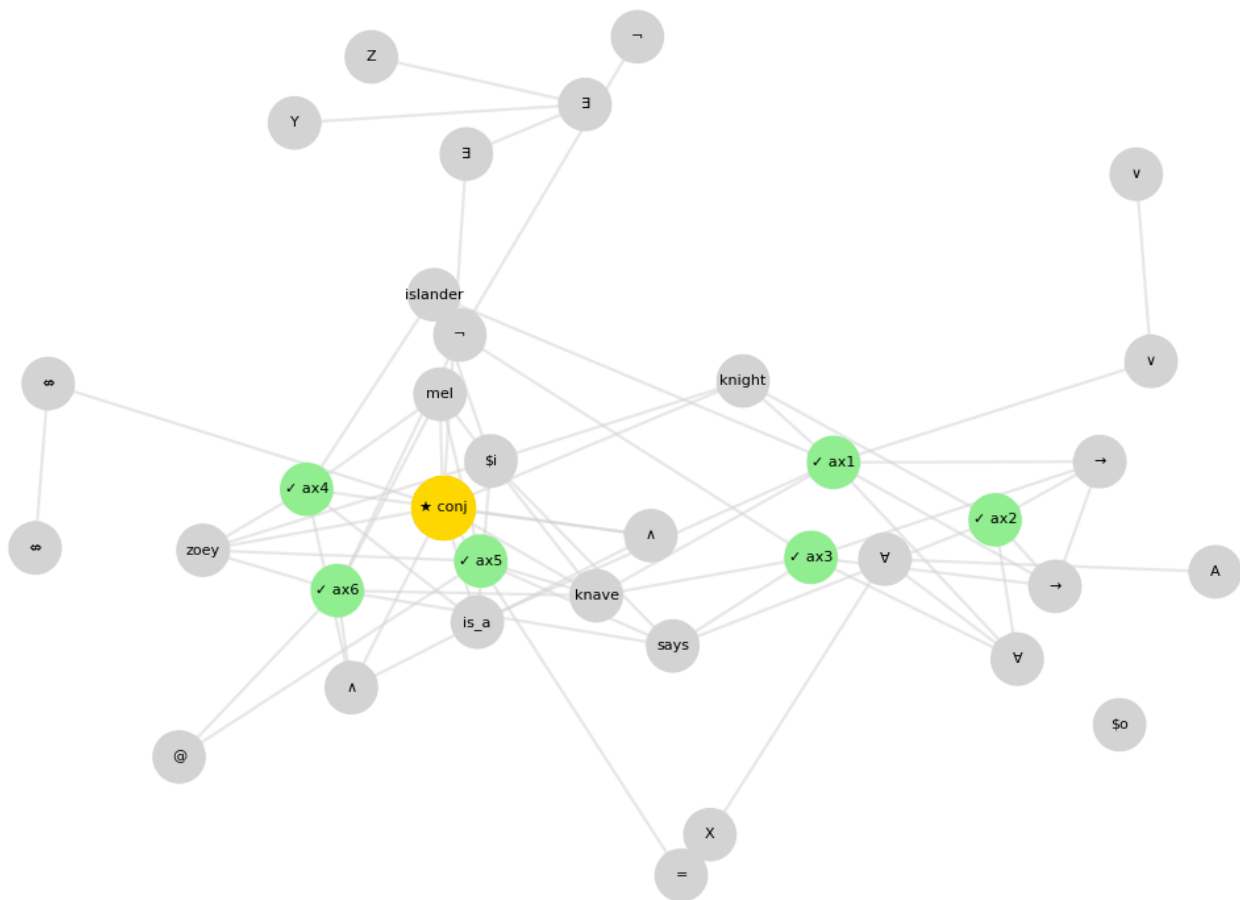
Edges

- CLAUSE (axiom)
- CLAUSE (conjecture)

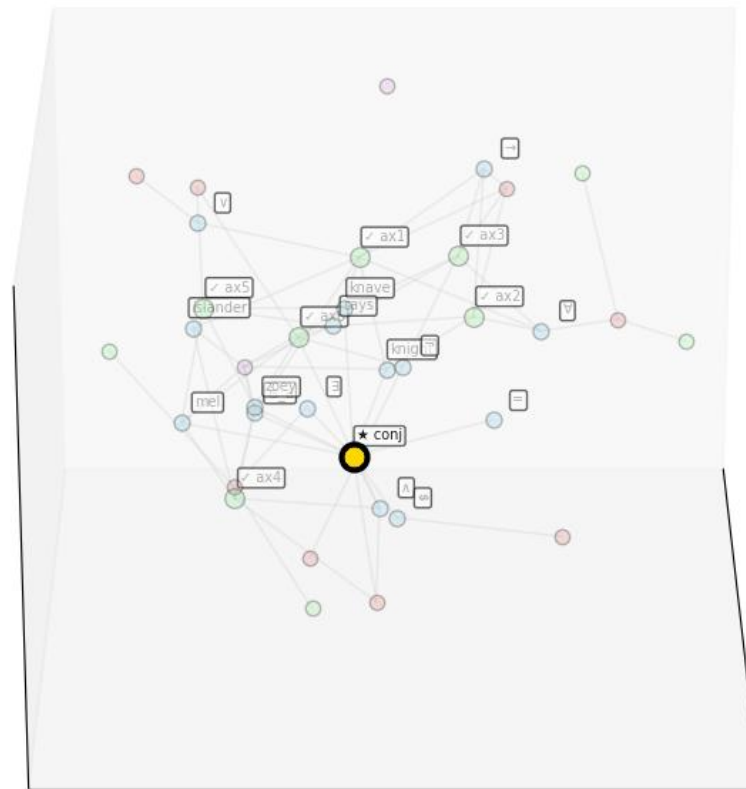
## Building Heterogeneous Graph for PUZ081^1 Step 1: CLAUSE nodes



GNN-v1: Message Propagation (PUZ081^1) — Layer 1/6  
Orange = active | Blue = passed | Reached: 1/36 nodes



GNN-v1: 3D Message Propagation — Layer 1/6  
Orange = active | Blue = passed | 1/36 nodes | Rotating



# GNN\_v1

Load pt file  
from vampire

Build Graph

Initialization

Message  
Passing

Scoring

Loss

## Training

Graph of the training  
dataset

Fixed Trainig Graph

GNN Weights  
(random)

Training Embeddings

Loss & Gradients

## Testing

New Test Graph for  
the test set....

Fixed Graph Structure

Frozen weights

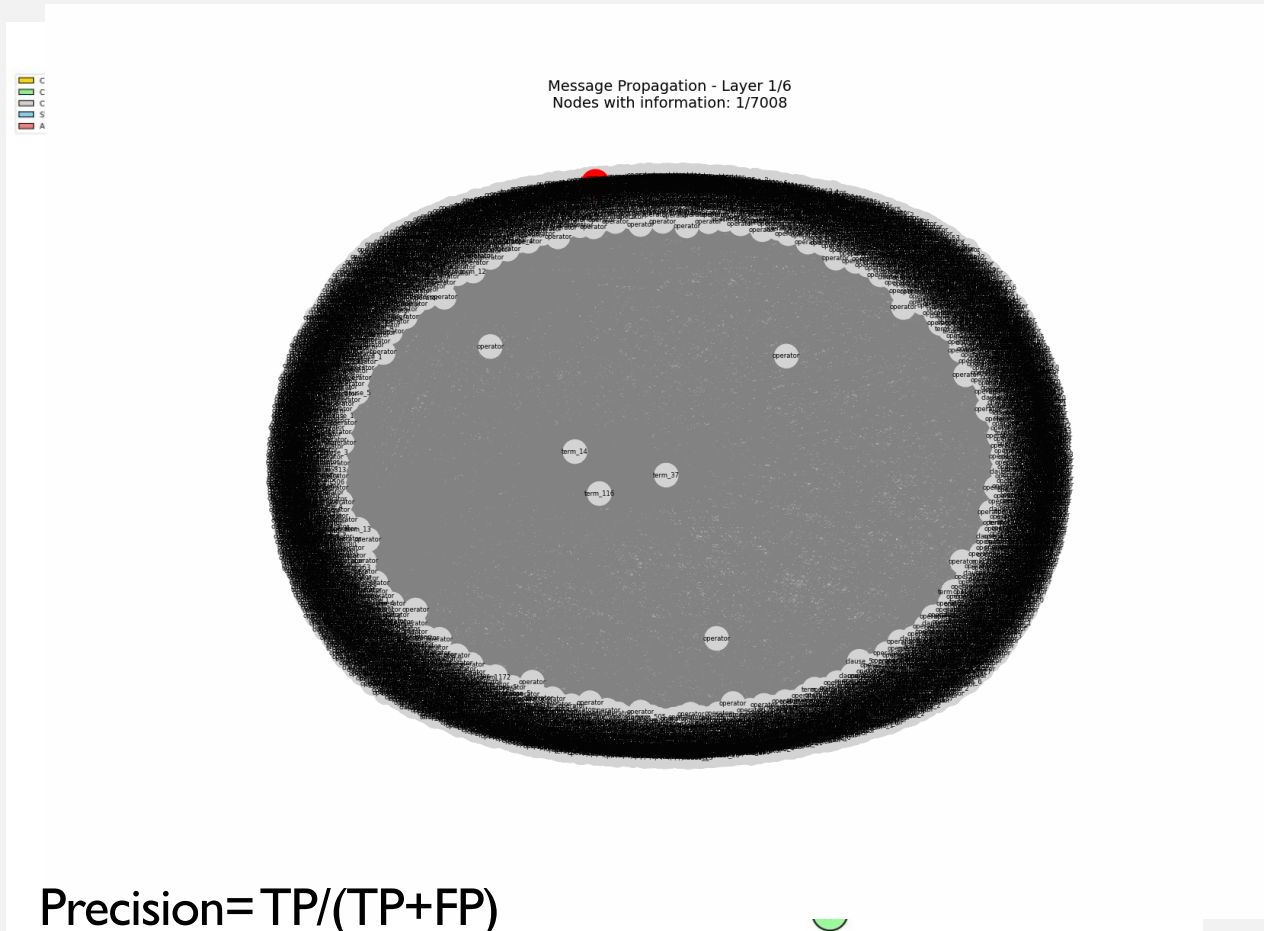
Inference Embeddings

Predictions

## Loss Function

1. Binary Cross Entropy Loss (BCE)
2. Focal Loss

Lin, Tsung-Yi, et al. "Focal loss for dense object detection." *Proceedings of the IEEE international conference on computer vision*. 2017.



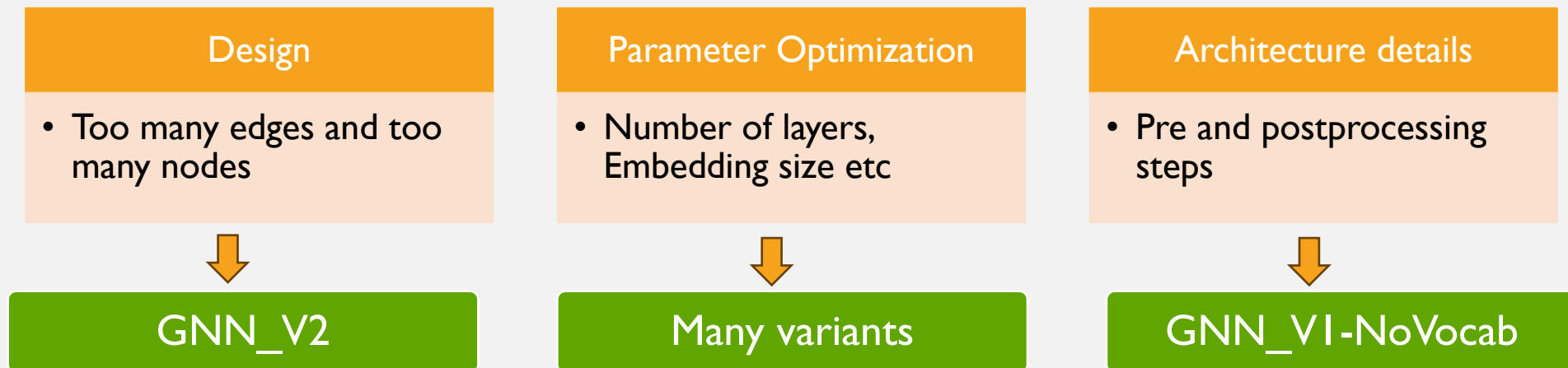
$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

$$\text{Accuracy} = \text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

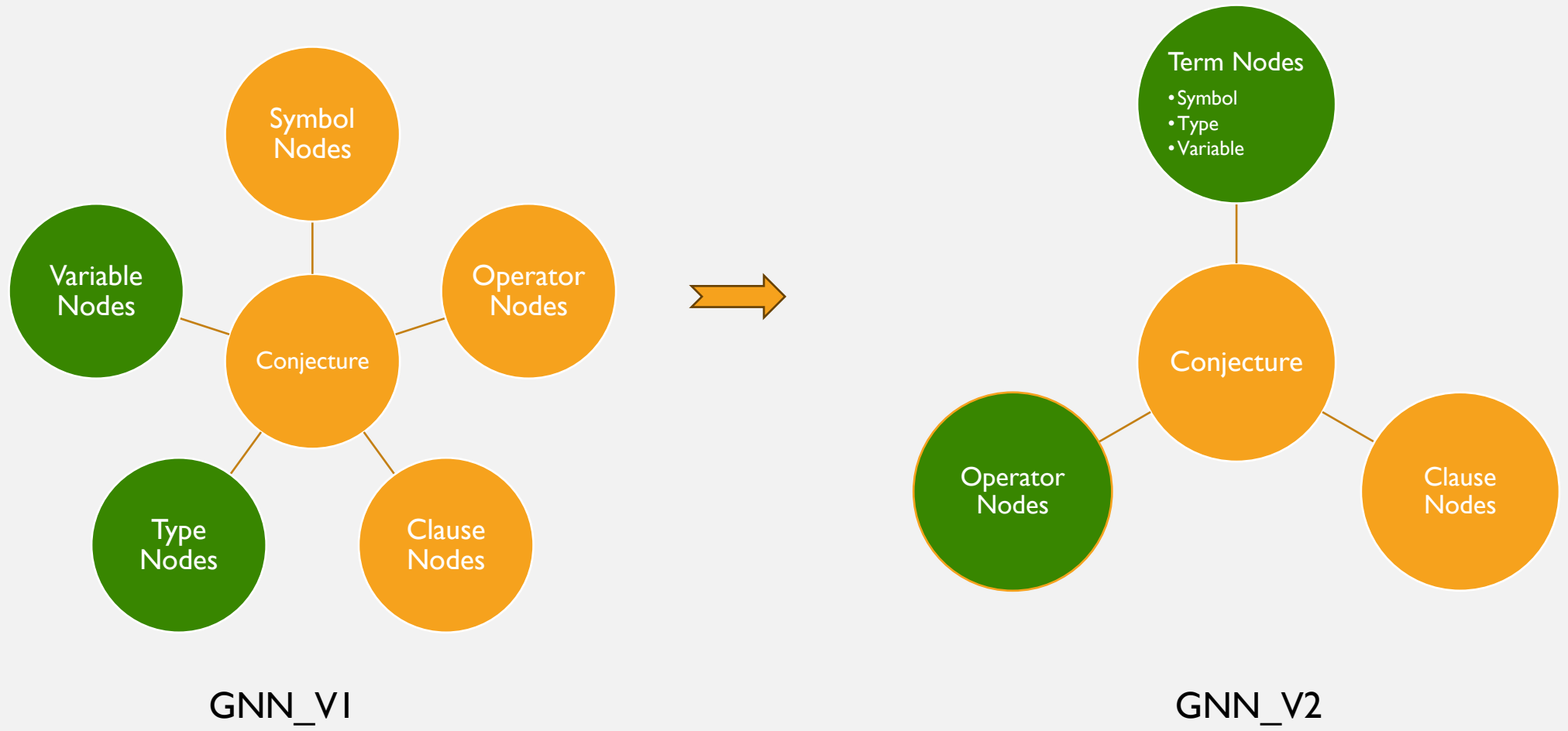
$$\text{F-score} = 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$$

# IMPROVING GNN-V1



Model	Input (.pt)	Representation / Preprocessing	Core Layers	Special Components	Output	Loss
GNN-v1	Graph data	5-type graph → Init	6 × SAGEConv	—	Score	Focal Loss
GNN-v1	Graph data	5-type graph → Init	6 × SAGEConv	GNN_V1_NoVocab	Score	Focal Loss
GNN-v2	Graph data	3-type graph + 5D edge features → Init	6 × GATv2Conv	Attention mechanism	Score	Focal Loss

# METHODOLOGY: GNN\_V2



# GNN\_v2

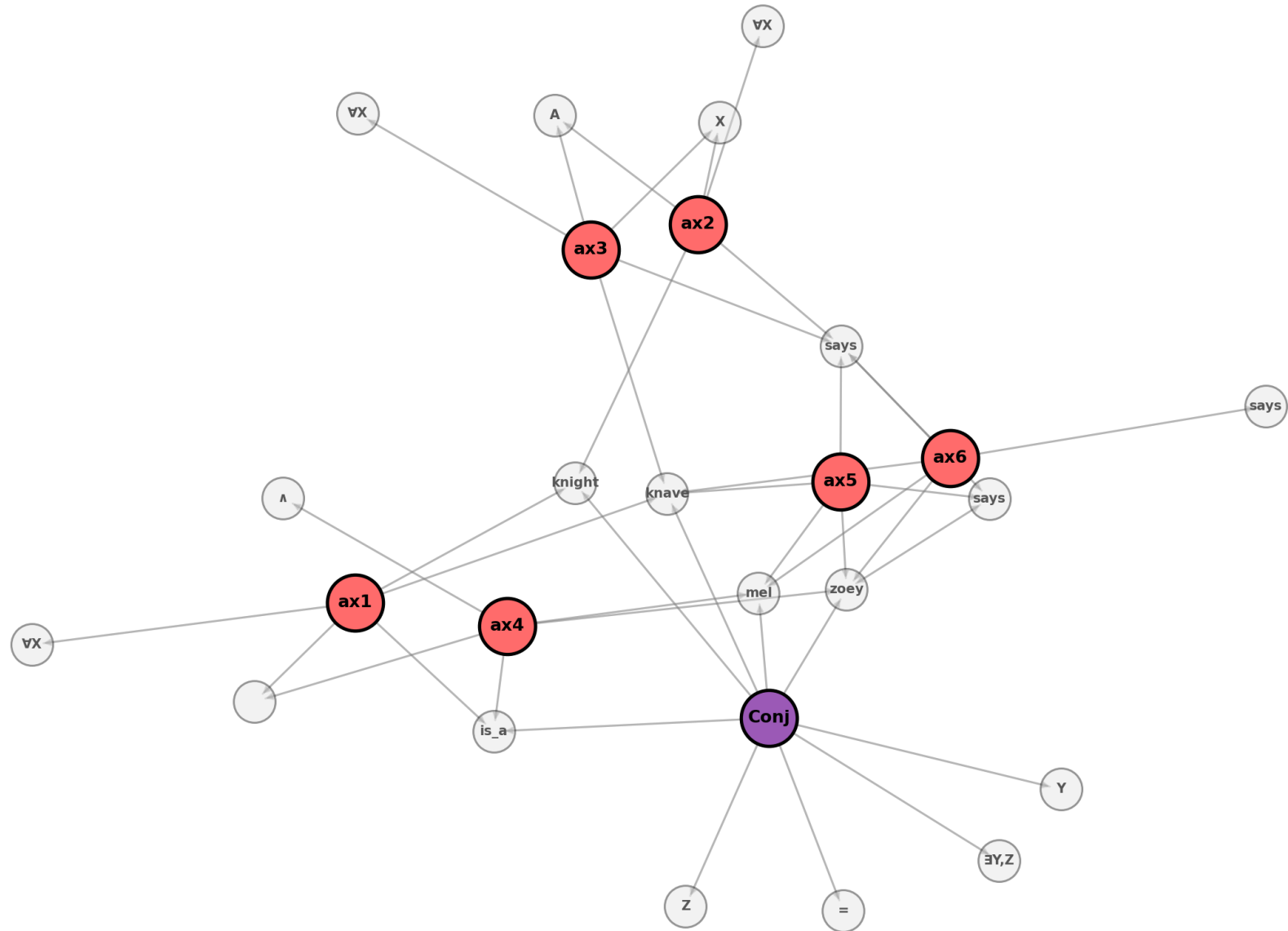
Clause nodes

Term nodes  
(symbol,  
variable, type)

Operator nodes

Edges with 5D  
edge Features

Node Type: CLAUSE (axioms & conjecture)



# GNN\_v2

Clause nodes



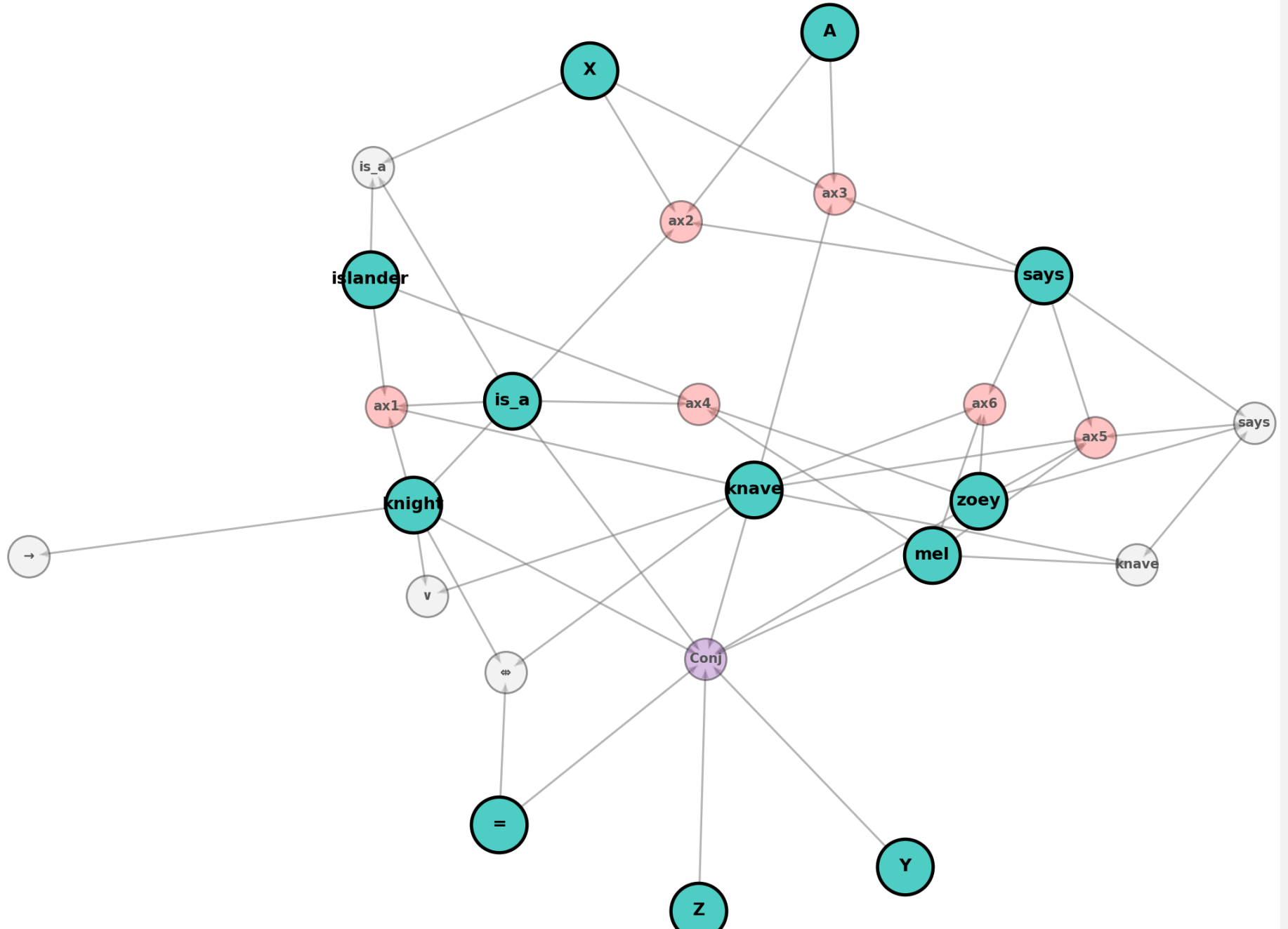
Term nodes  
(symbol,  
variable, type)



Operator nodes



Edges with 5D  
edge Features







# GNN\_V2

Term nodes  
(symbol,  
variable, type)



Operator  
nodes



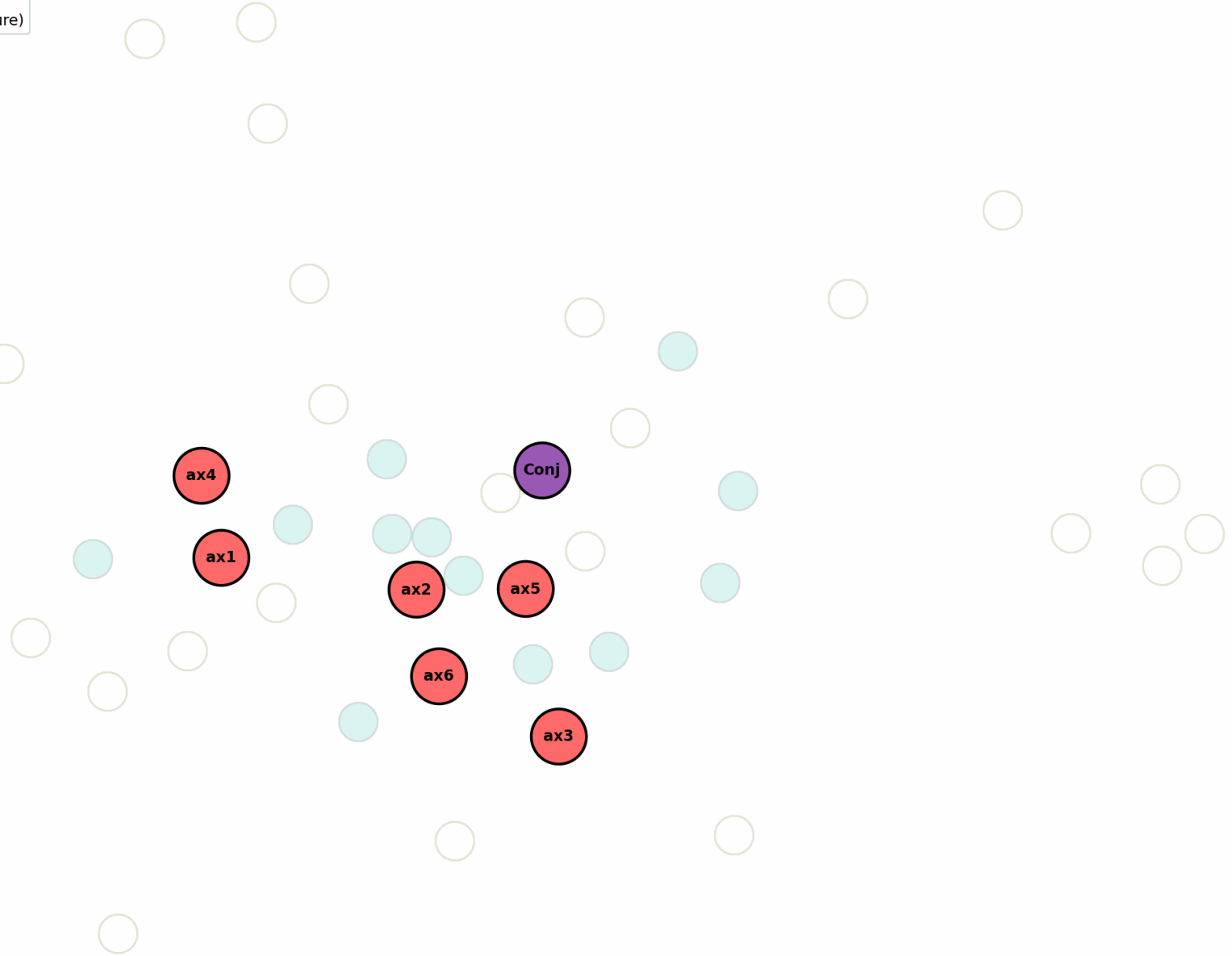
Clause nodes



Edges with 5D  
edge Features

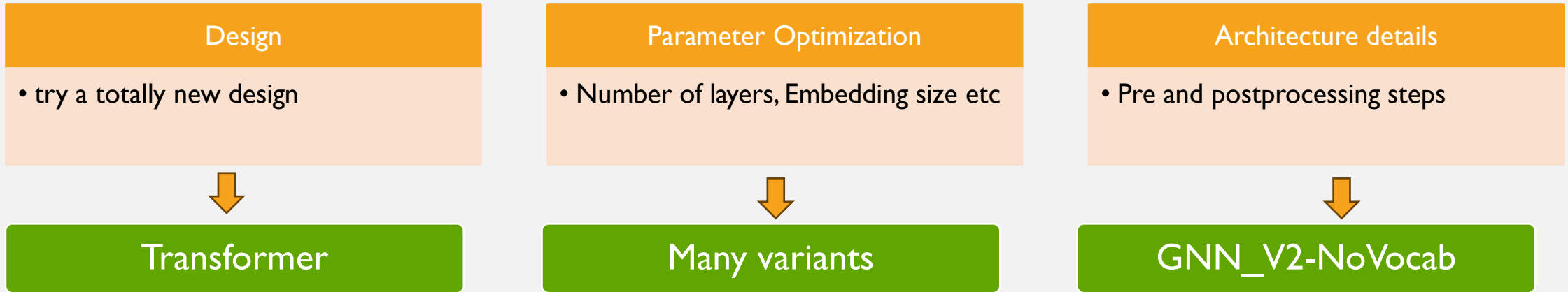
## GNN v2 Graph Construction Step 1: CLAUSE nodes

- CLAUSE (axiom)
- CLAUSE (conjecture)



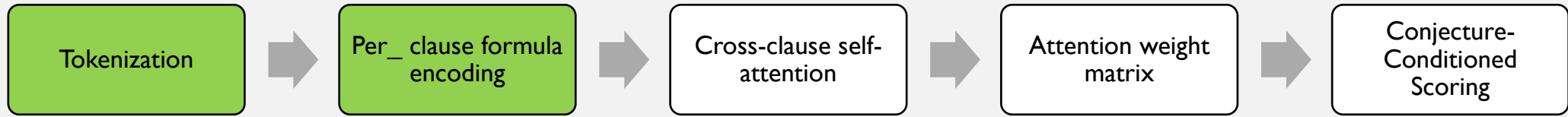


# Improving GNN\_v2

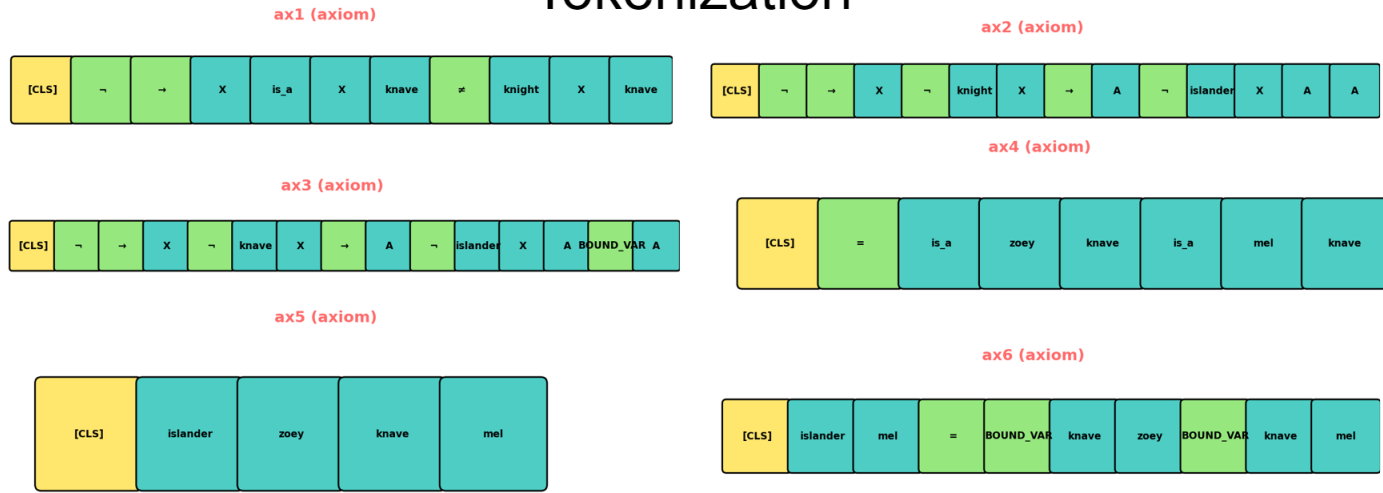


Model	Input (.pt)	Representation / Preprocessing	Core Layers	Special Components	Output	Loss
GNN_V1	Graph data	5-type graph → Init	6 × SAGEConv	—	Score	Focal Loss
GNN_V1_NoVocab	Graph data	5-type graph → Init	6 × SAGEConv	GNN_V1_NoVocab	Score	Focal Loss
GNN_V2	Graph data	3-type graph + 5D edge features → Init	6 × GATv2Conv	Attention mechanism	Score	Focal Loss
GNN_V2_NoVocab	Graph data	3-type graph + 5D edge features → Init	6 × GATv2Conv	Attention mechanism	Score	Focal Loss

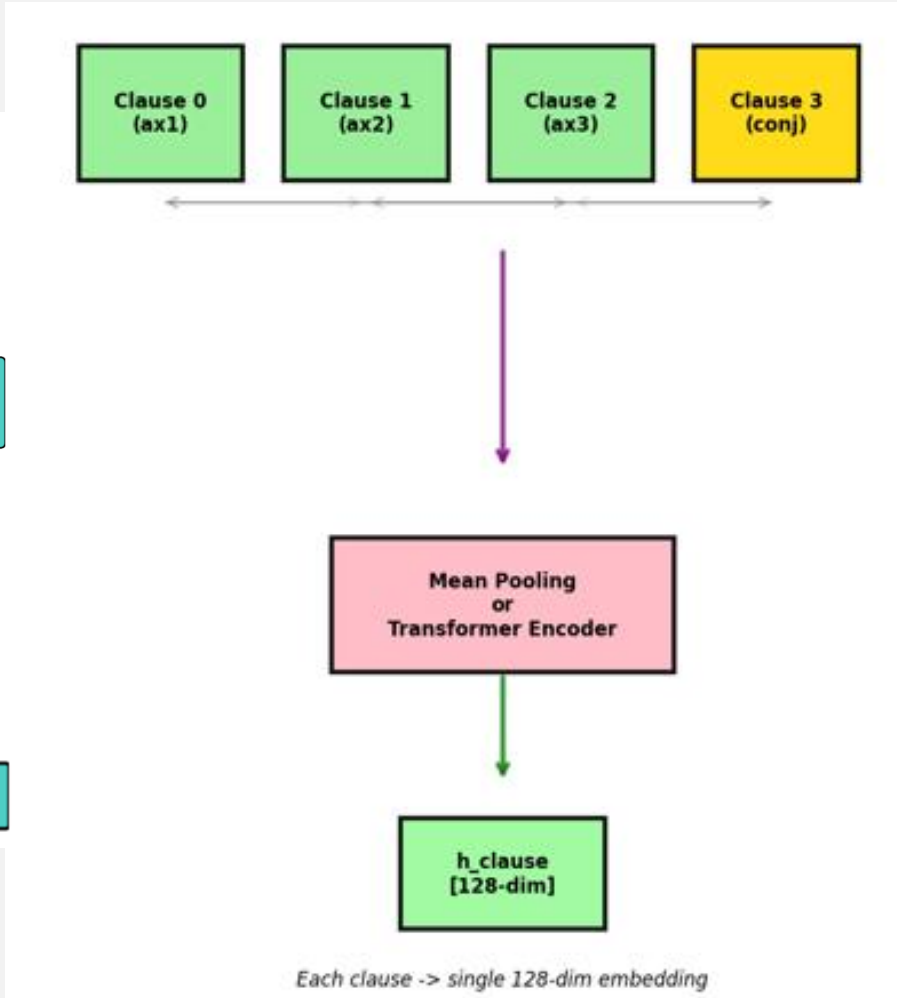
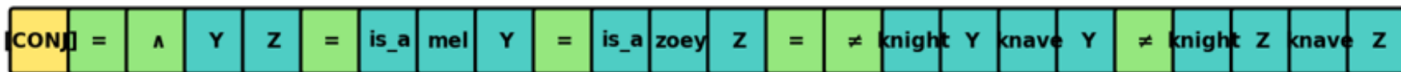
# Transformer



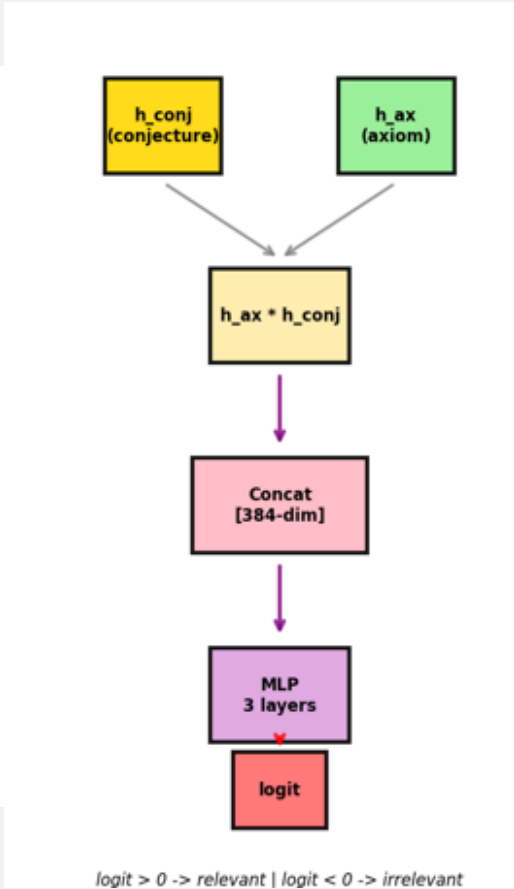
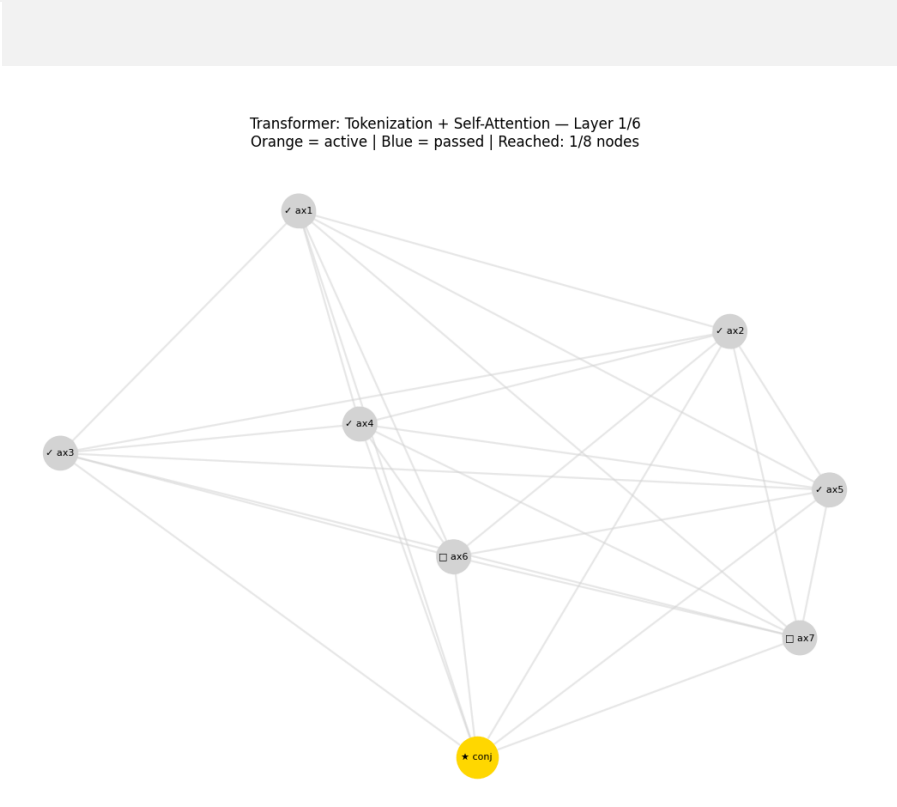
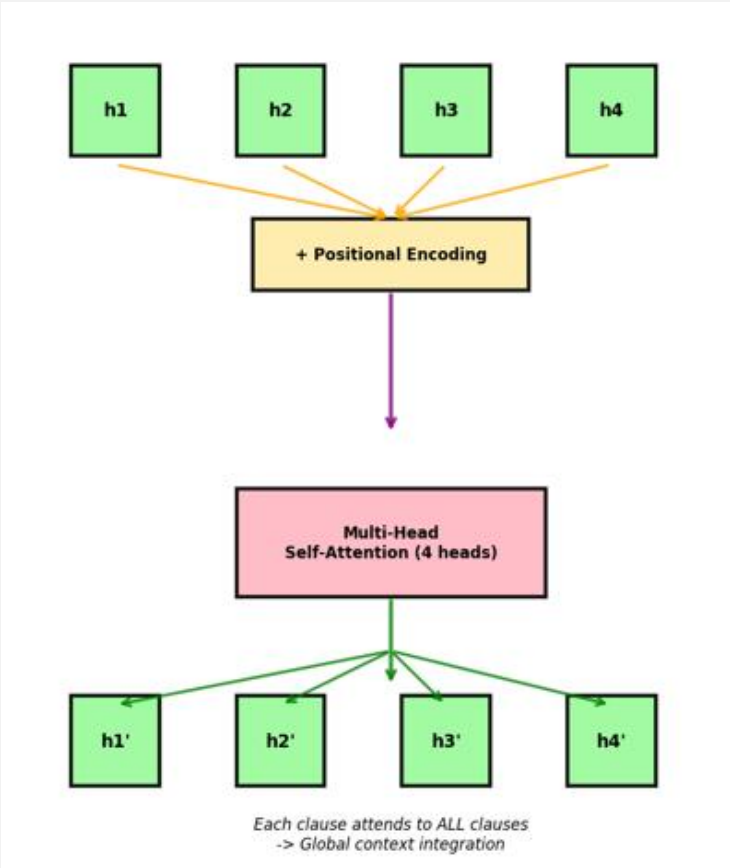
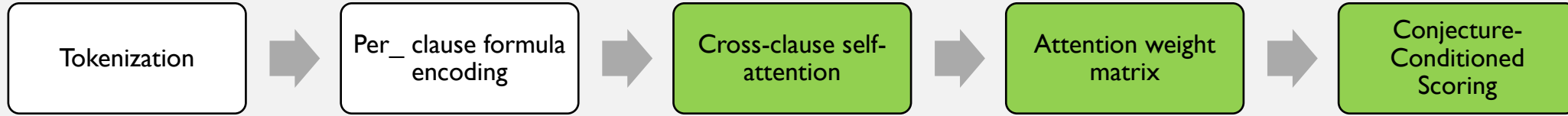
## Tokenization



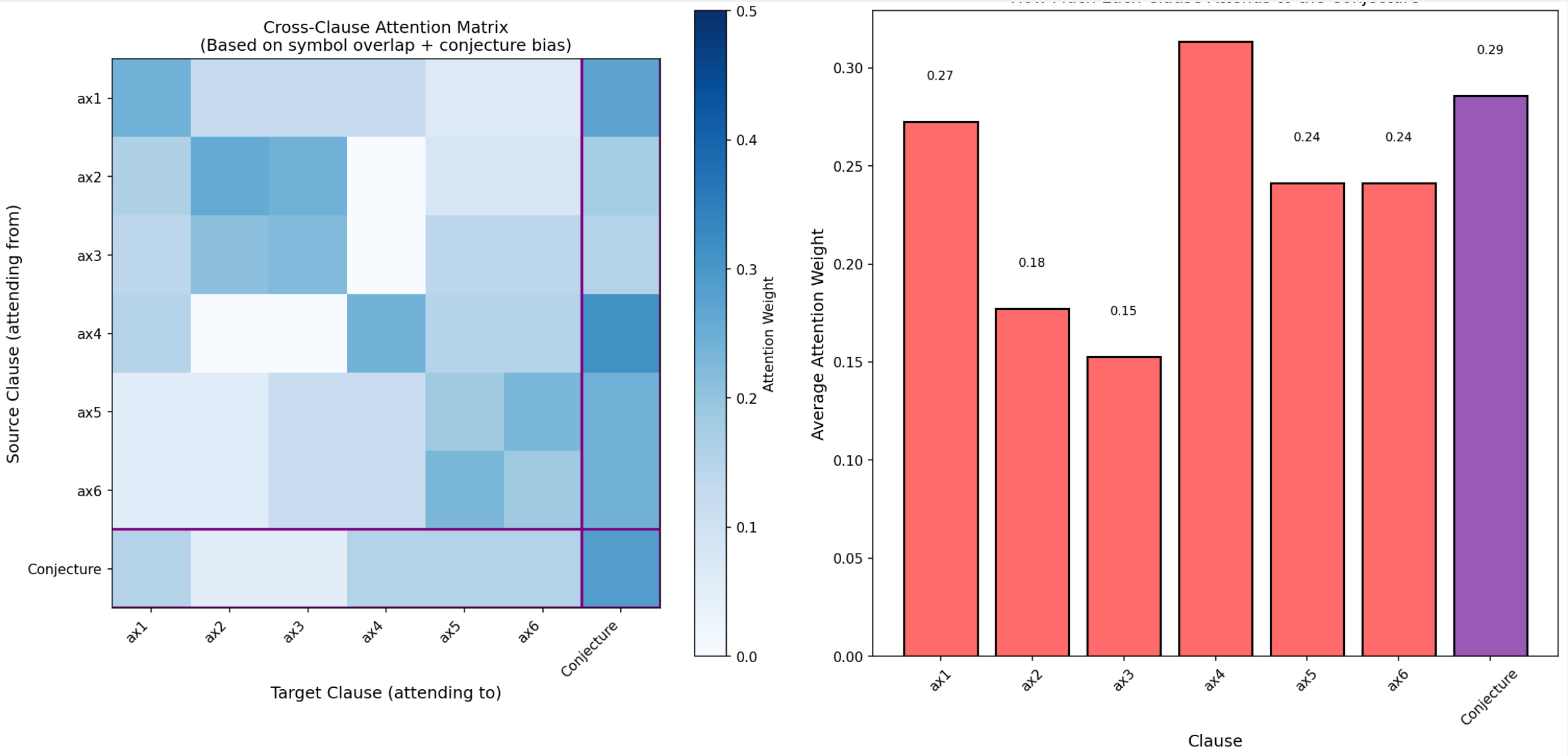
## Conjecture



# Transformer



# Transformer: Self-Attention Weight Matrix



# Improve it?

# IMPROVE IT?

## Design

- Type information is lost



Contrastive Transformer

## Parameter Optimization

- Number of layers, Embedding size etc



Many variants

## Architecture details

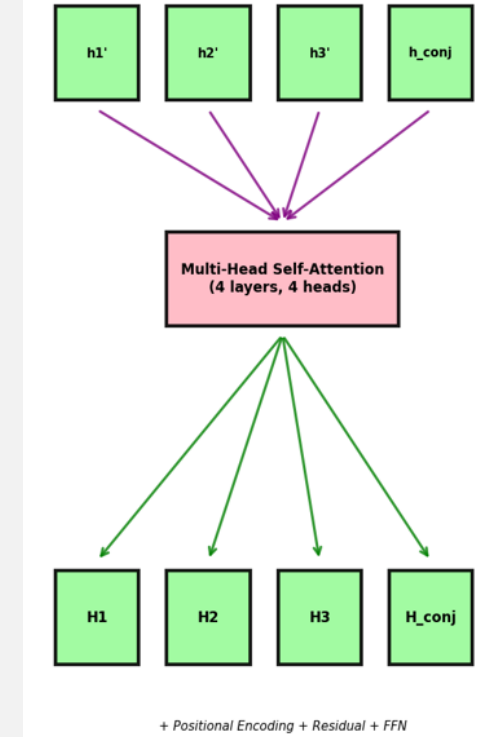
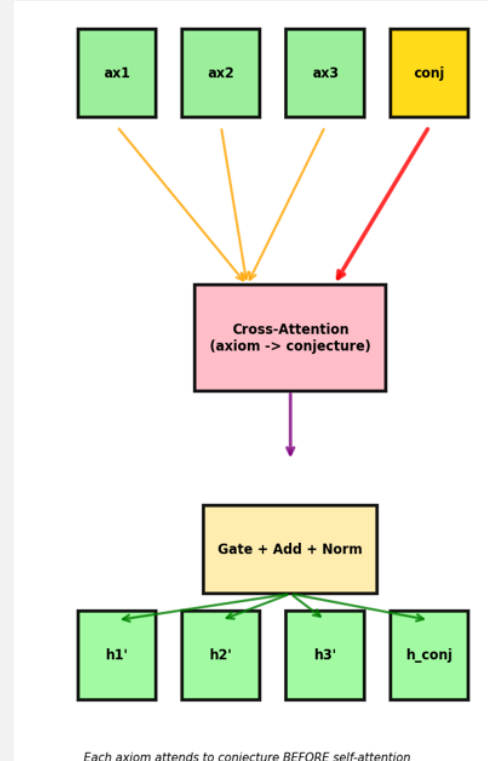
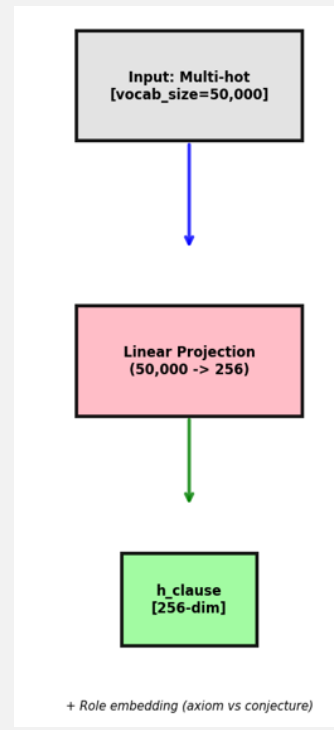
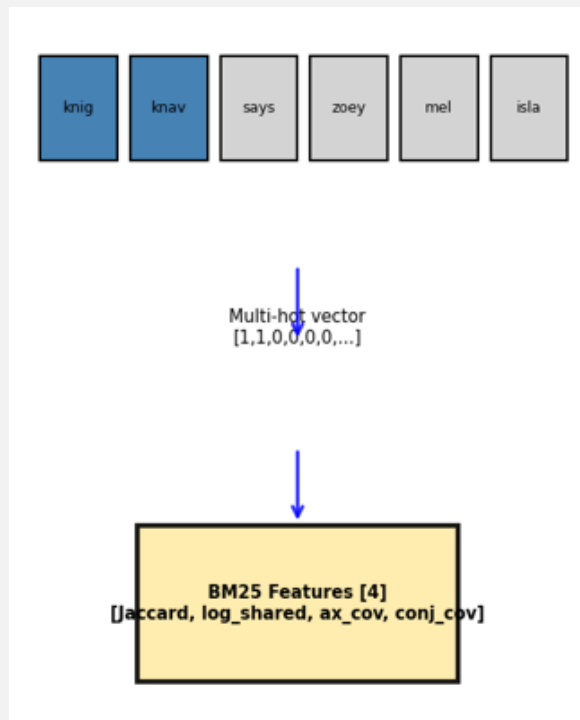
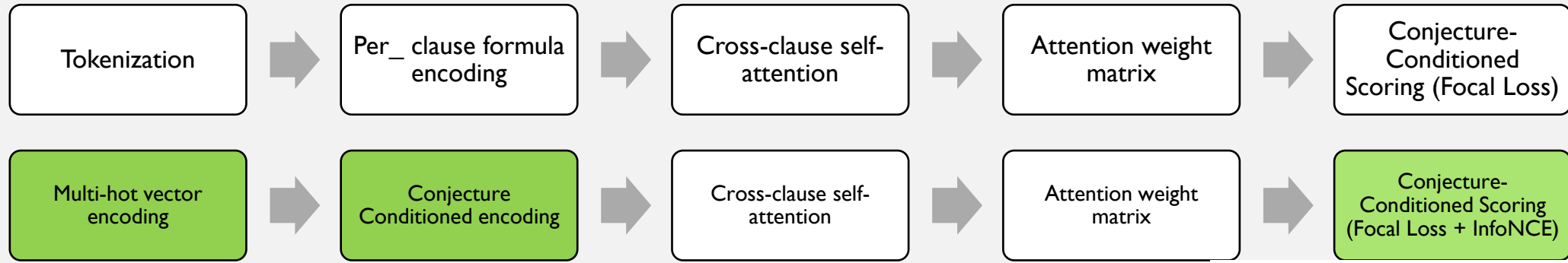
- Pre and postprocessing steps



w/o Positional Encoding,  
Typed Transformer

Model	Input (.pt)	Representation / Preprocessing	Core Layers	Special Components	Output	Loss
Transformer	Token data	Linearize → Token Embedding	Formula Encoder + 4 × Cross-Attention	Pooling / Transformer encoder	Score	—
Typed_ Transformer	Token data with type info	Linearize → Token Embedding	Formula Encoder + 4 × Cross-Attention	Pooling / Transformer encoder	Score	—
Transformer_w/o_ Positional Encoding	Token data	Linearize → Token Embedding	Formula Encoder w/o Positional encoding + 4 × Cross-Attention	Pooling / Transformer encoder	Score	—
Contrastive	Feature data	Multi-hot + BM25 (4) → Projection	2 × Cross-Attention + 4 × Self-Attention	Axiom → Conjecture alignment	Dual Score	Focal Loss + InfoNCE

# Contrastive Transformer



Mikuła, Maciej, et al. "Magnushammer: A transformer-based approach to premise selection." *International Conference on Learning Representations*. Vol. 2024. 2024.

# Contrastive Transformer

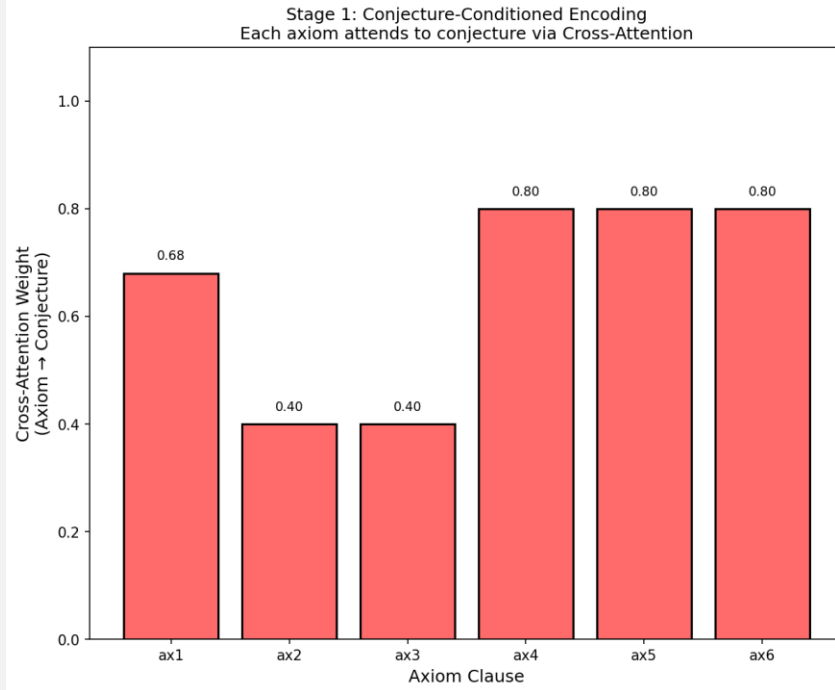
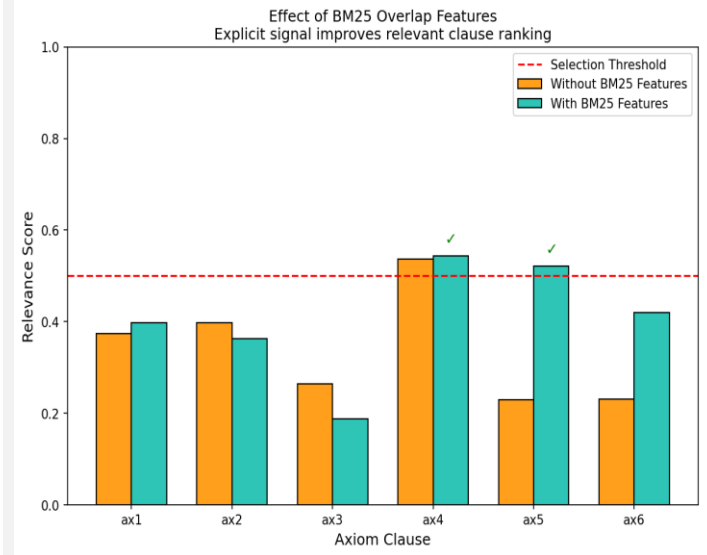
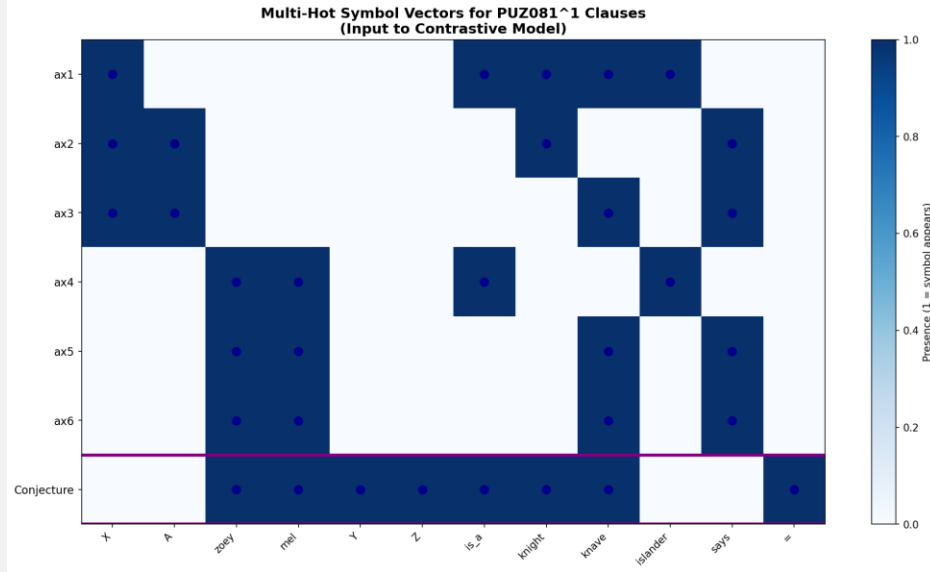
Multi-hot vector encoding

Conjecture Conditioned encoding

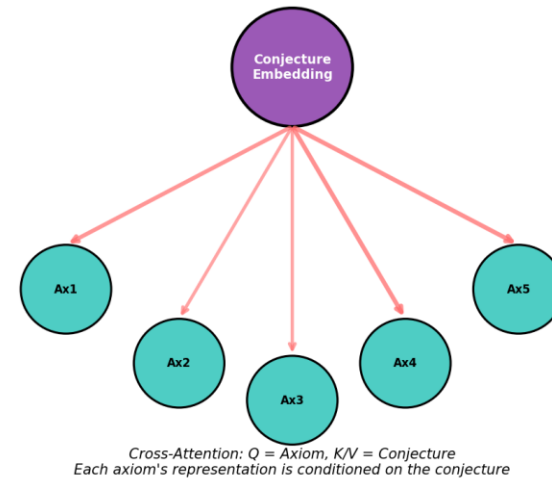
Cross-clause self-attention

Attention weight matrix

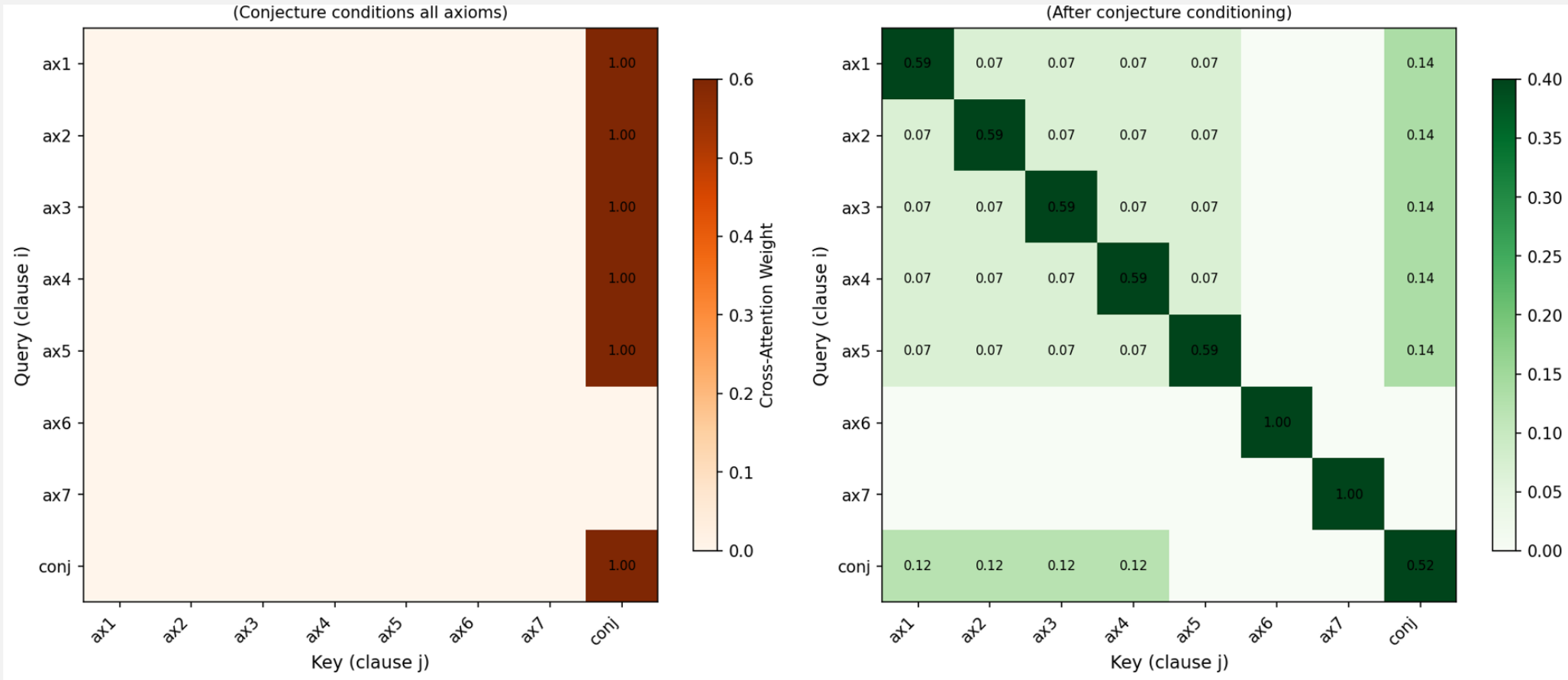
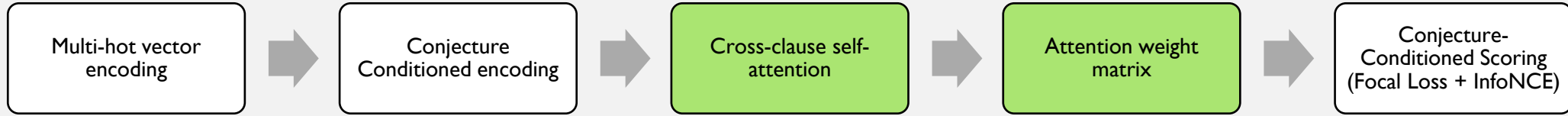
Conjecture-Conditioned Scoring (Focal Loss + InfoNCE)



Conjecture-Conditioned Encoding (ARCG-NN style)

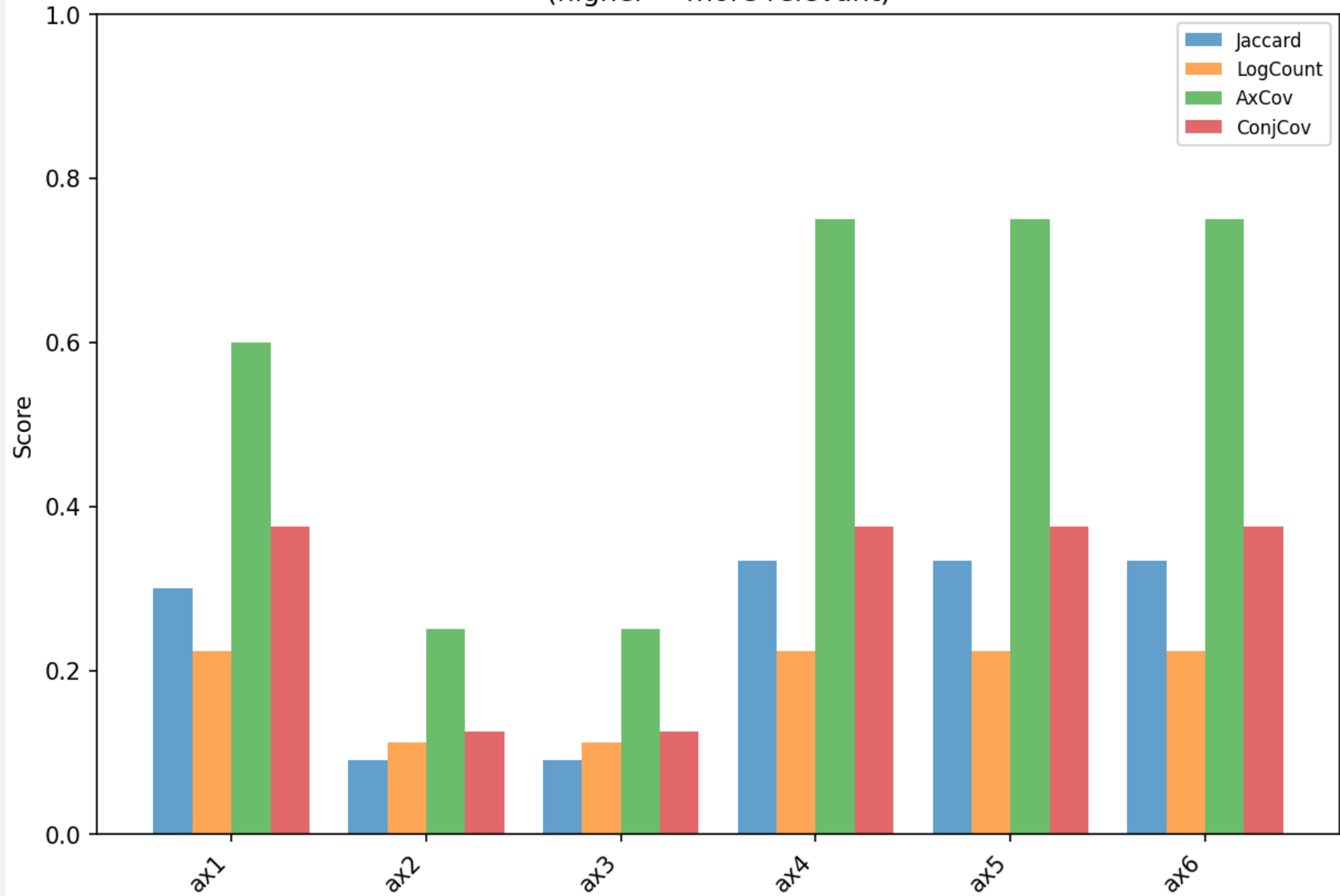


# Contrastive Transformer



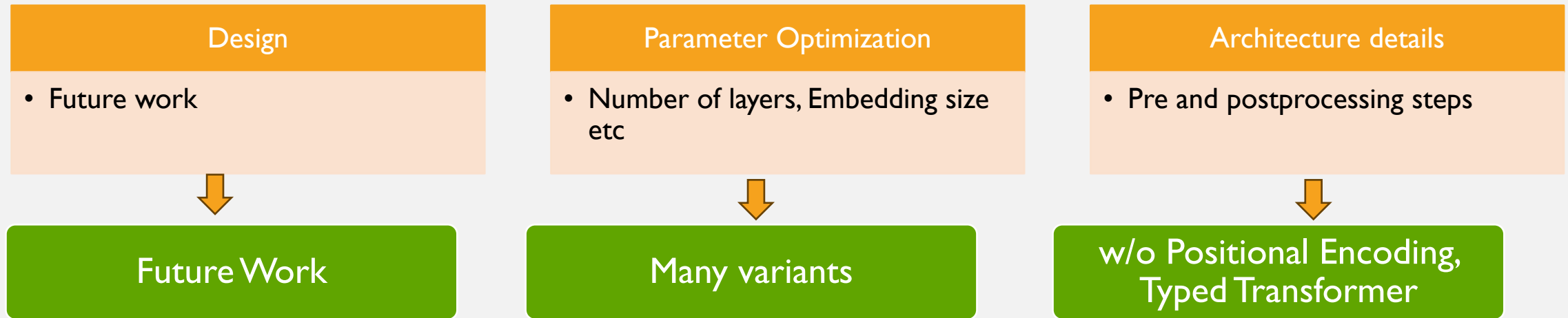
CROSS ATTENTION: AXIOM TO CONJECTURE  
 SELF ATTENTION ON ALL AXIOMS

BM25 Overlap Features per Axiom  
(higher = more relevant)



# Improve it?

# IMPROVE IT?



Model	Input (.pt)	Representation / Preprocessing	Core Layers	Special Components	Output	Loss
Contrastive	Feature data	Multi-hot + BM25 (4) → Projection	2 × Cross-Attention + 4 × Self-Attention	Axiom → Conjecture alignment	Dual Score	Focal Loss + InfoNCE
Typed_Contrastive	Feature data including type	Multi-hot + BM25 (4) → Projection	2 × Cross-Attention + 4 × Self-Attention	Axiom → Conjecture alignment	Dual Score	Focal Loss + InfoNCE
Contrastive_No Positional Encoding	Feature data	Multi-hot + BM25 (4) → Projection	2 × Cross-Attention + 4 × Self-Attention, No Positional Encoding	Axiom → Conjecture alignment	Dual Score	Focal Loss + InfoNCE

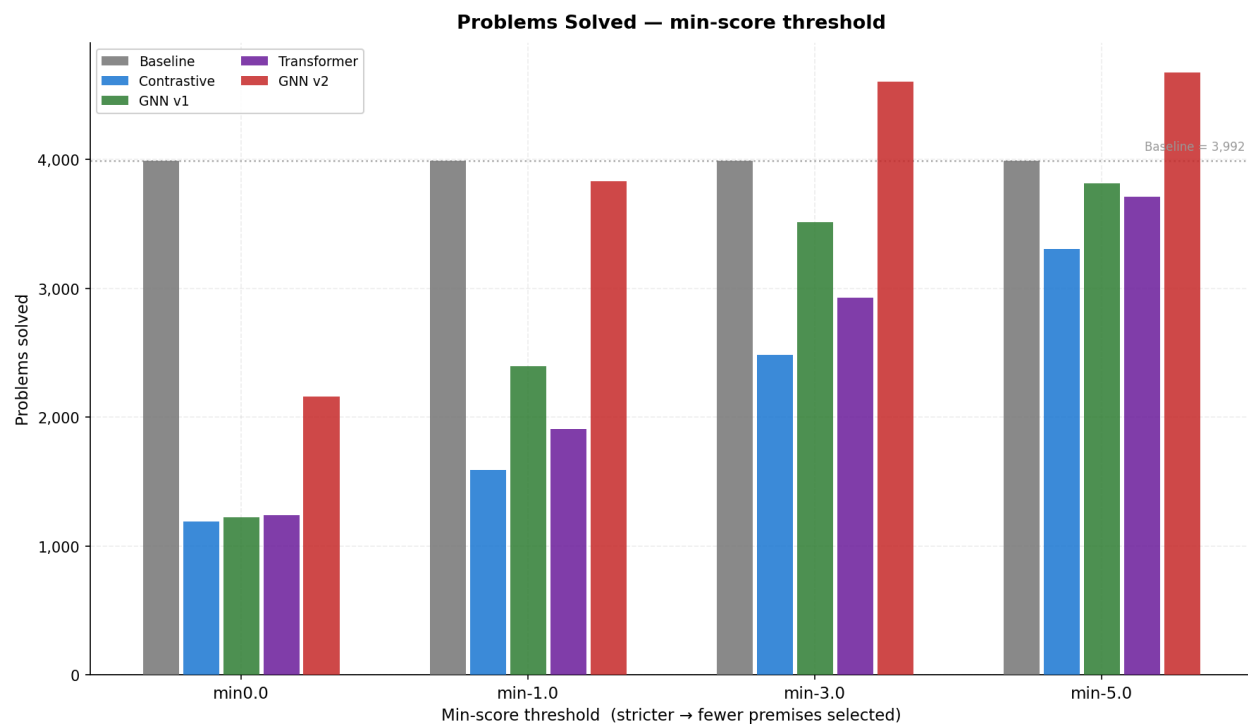
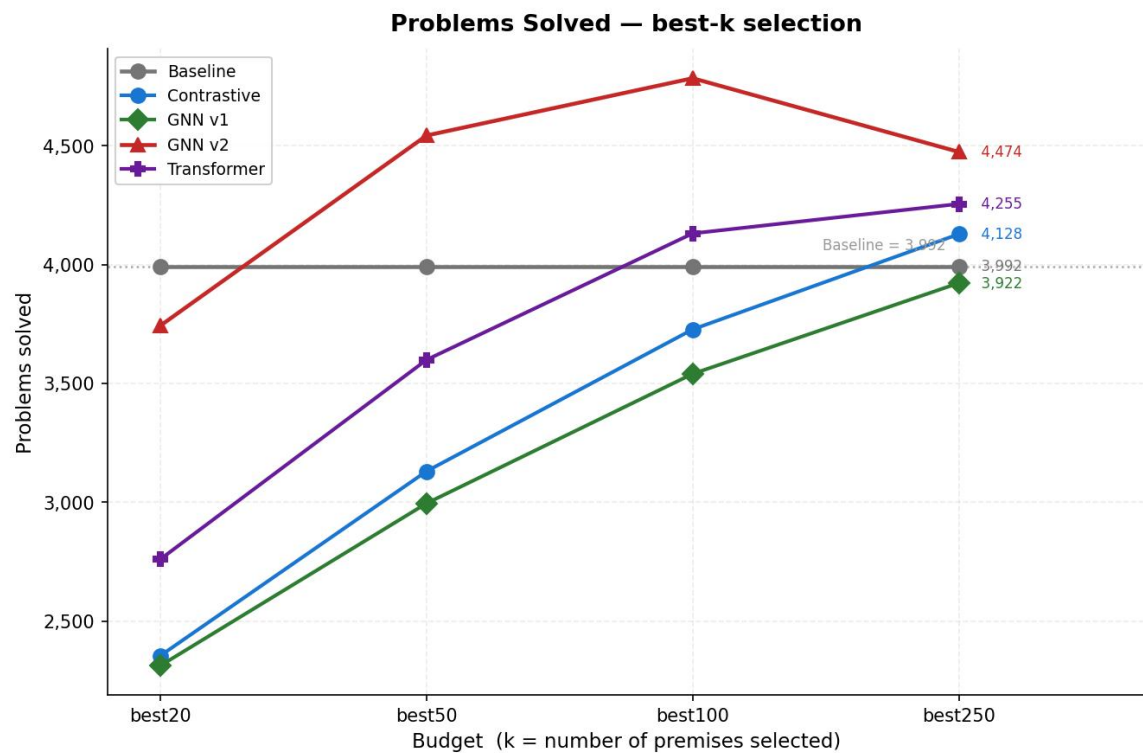
# EXPERIMENTS

- Two strategies:
  - Best X axioms (best 10,20, 50,100, 250, 512(baseline))
  - All axioms above Logit Score X (min 0.5, 0.0 , -0.5, -1.0 , -2.0, -3.0, -5.0)
  - T-Time (1s, 10s)
- Dataset:
  - Training: 26,000 problems from Isabelle (80/20 train/val split)
  - Testing : 8900 holdout problems

## Model

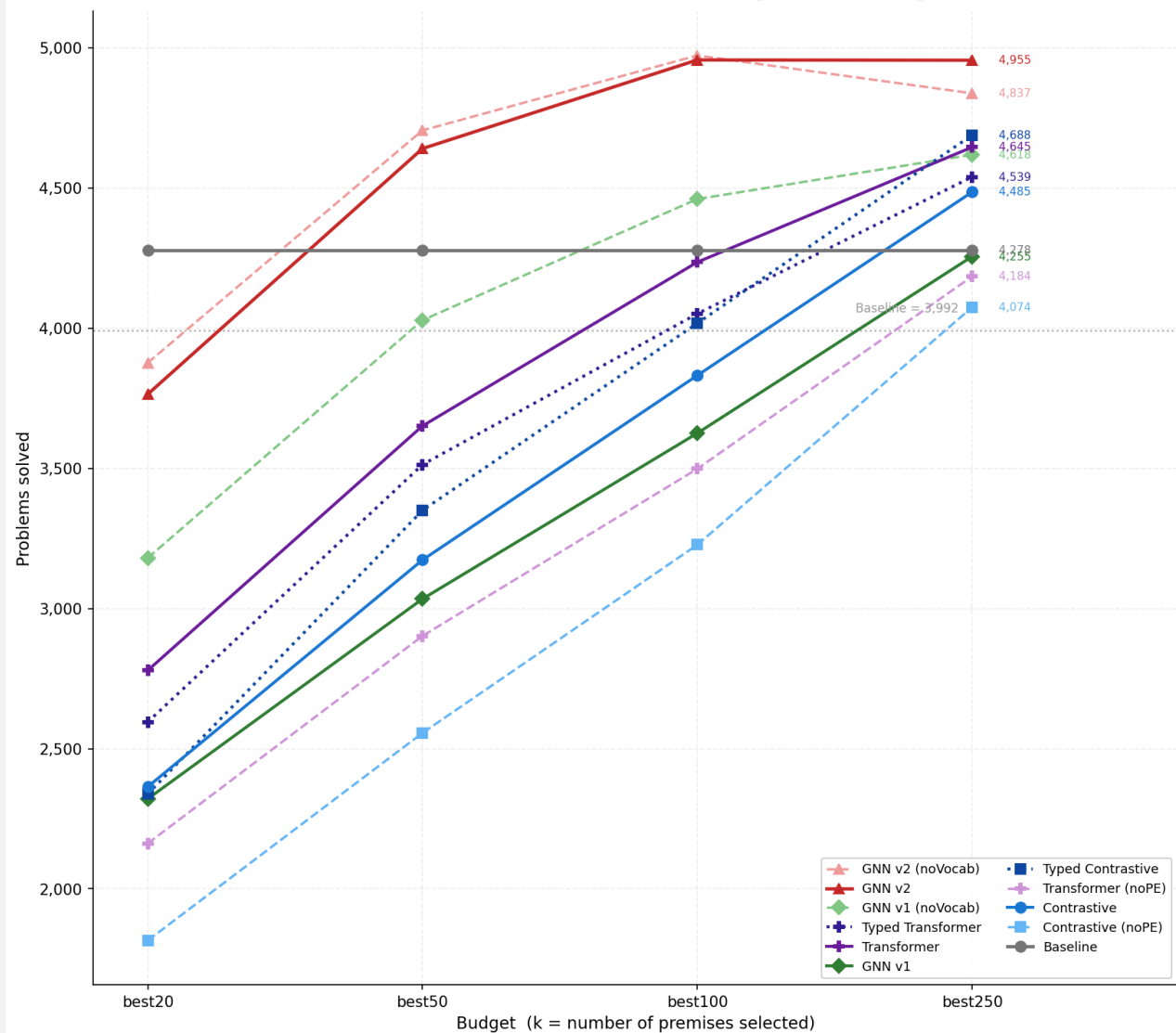
- Baseline (MePo)
  - GNN\_V1
  - GNN\_V1\_NoVocab
- GNN-V2
  - GNN\_V2\_NoVocab
- Transformer
  - TypedTransformer
  - Transformer\_NoPE
- Contrastive
  - Typed Contrastive Transformer
  - Contrastive Transformer NoPE

# INITIAL RESULTS (T=1S)

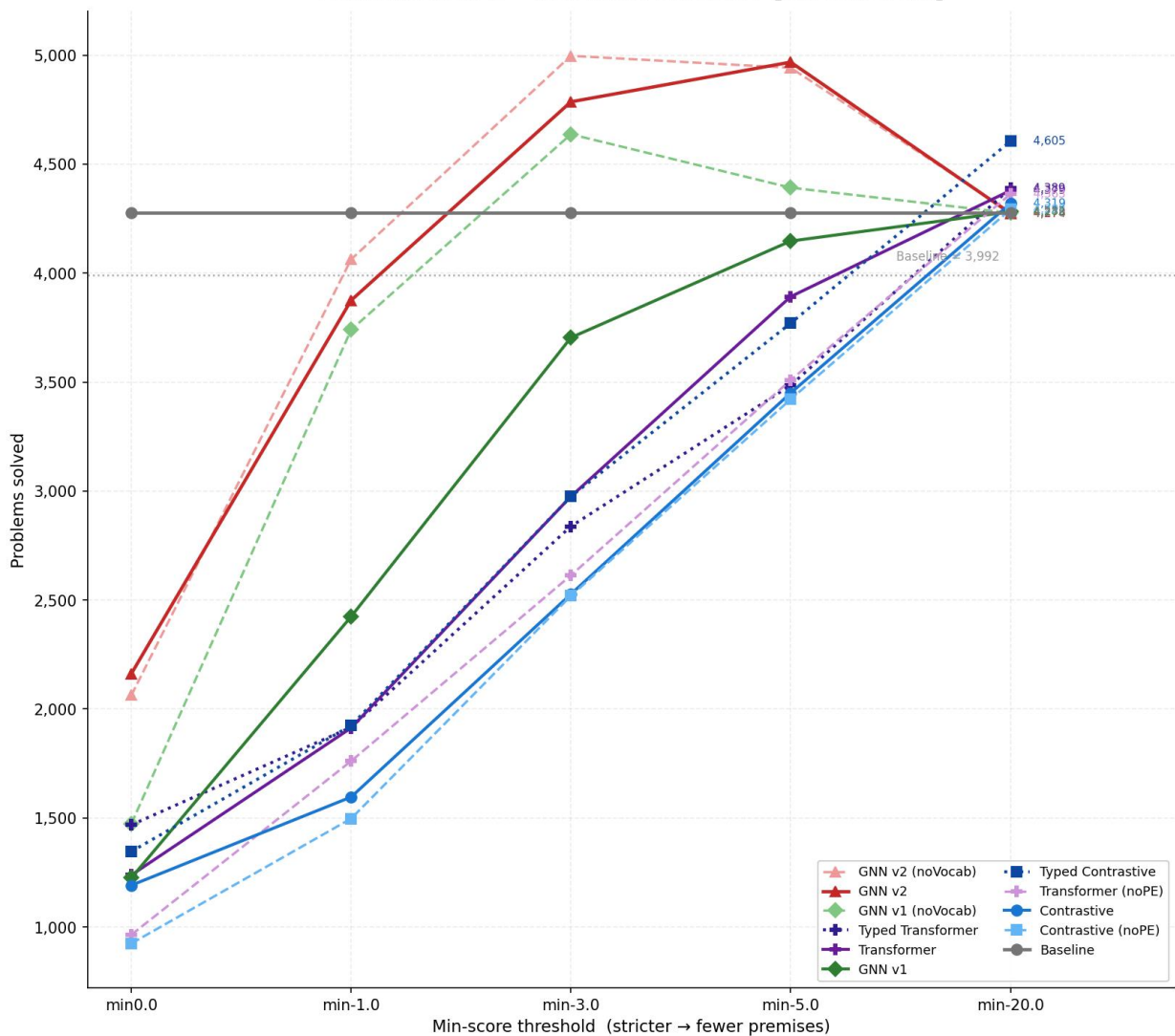


# INITIAL RESULTS (T=10 S)

Problems Solved — best-k selection [full model set]

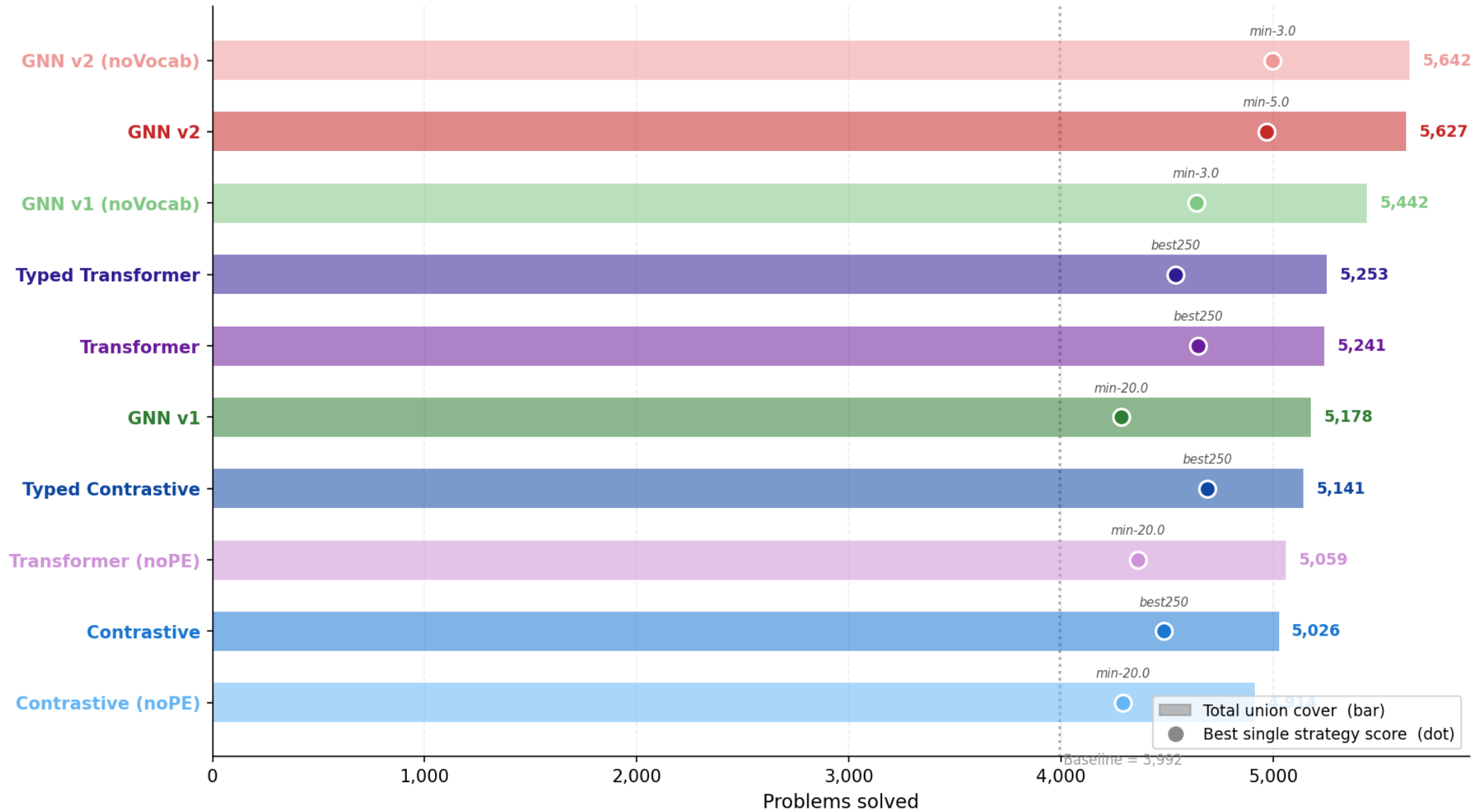


Problems Solved — min-score threshold [full model set]

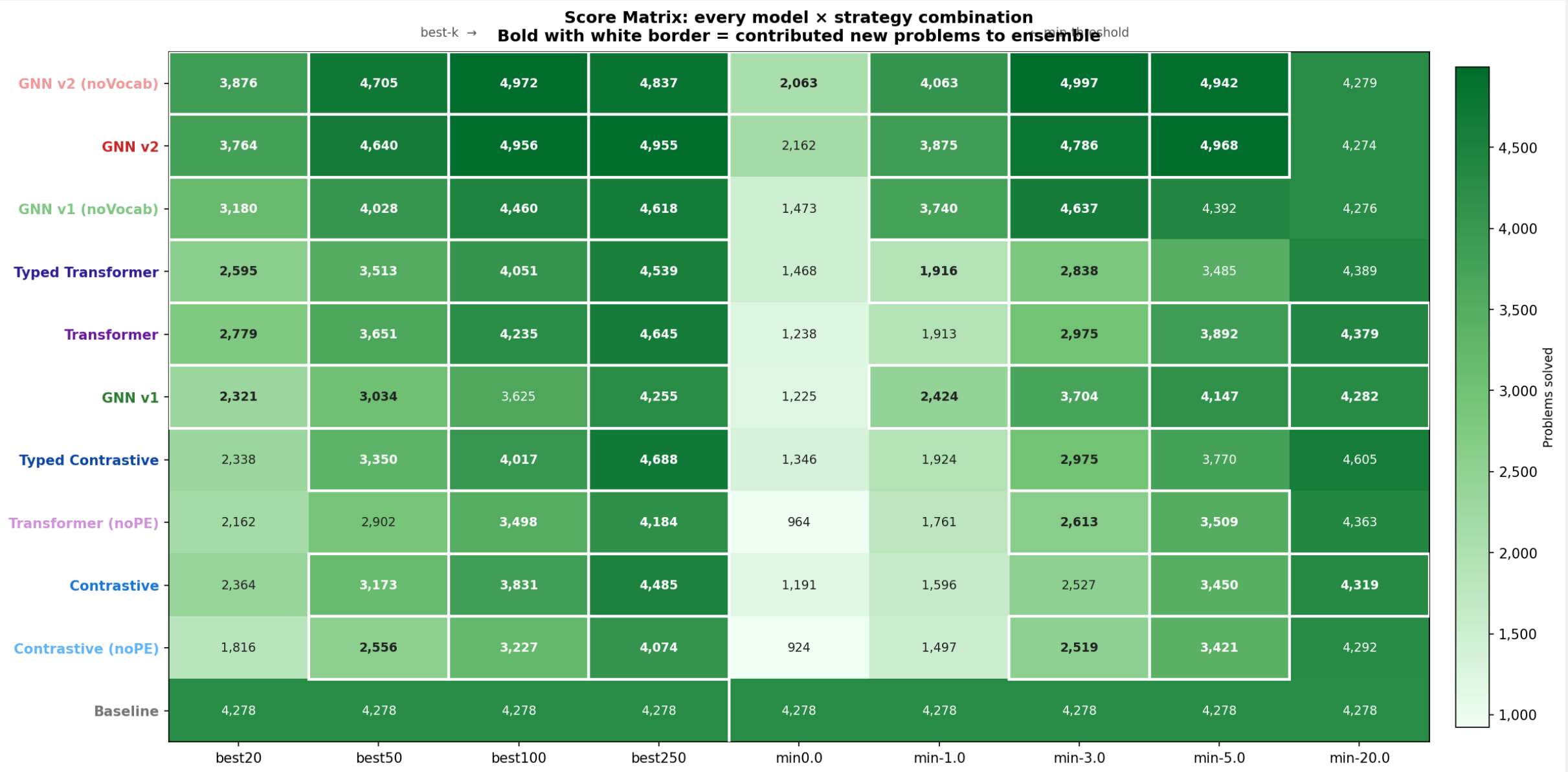


# INITIAL RESULTS (T=10 S)

## Total Union Coverage + Best Single Strategy Score



# RESULTS (T=10 S, GREEDY COVER)



RESULTS (T=10 S, GREEDY COVER)

MODEL	VARIANT	SCORE	+NEW	+%	TOTAL	COVER%
gnn_v2_novocab	min-3.0	4997	4997	—	4997	55.65%
gnn_v2_novocab	best250	4837	333	6.66%	5330	59.35%
gnn_v2	best50	4640	158	2.96%	5488	61.11%
transformer	best250	4645	78	1.42%	5566	61.98%
gnn_v2	best100	4956	57	1.02%	5623	62.62%
gnn_v2_novocab	min-1.0	4063	45	0.80%	5668	63.12%
gnn_v1_novocab	best100	4460	32	0.56%	5700	63.47%
gnn_v1	min-20.0	4282	26	0.46%	5726	63.76%
transformer	best50	3651	21	0.37%	5747	64.00%
gnn_v2	min-5.0	4968	19	0.33%	5766	64.21%

# CONCLUSION & FUTURE WORK

Premise Selection is a Promising direction towards solving more HOL problems faster.

Our models are also promising but they could be better.

