

Adventures in Neurosymbolic Machine Learning

symbols to vectors
vectors to symbols
plus other stuff

Charles Sutton

9 April 2018

Companion Site: bit.ly/adventures-neurosymbolic

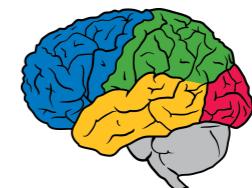


THE UNIVERSITY OF EDINBURGH
informatics

The
Alan Turing
Institute

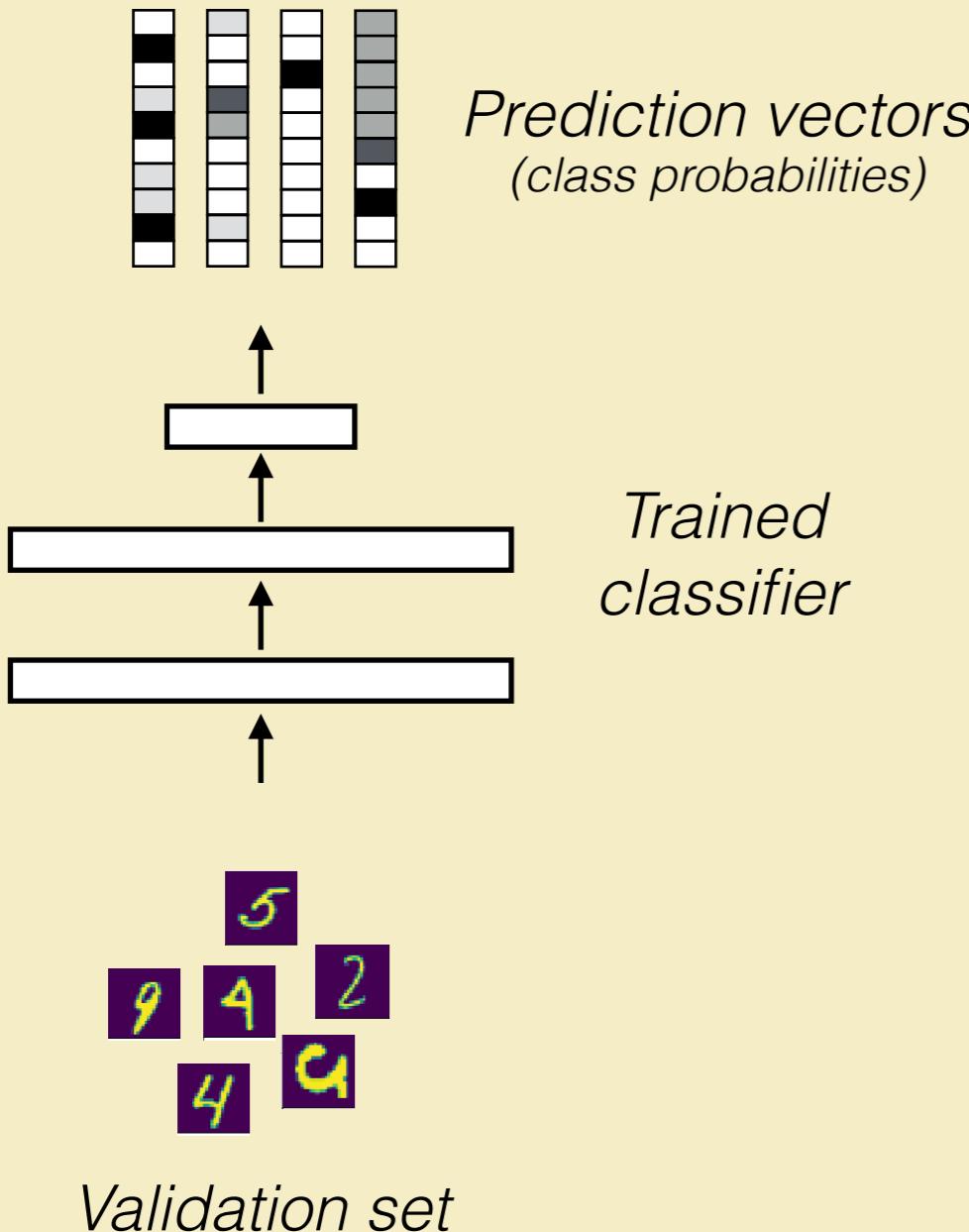
Microsoft
Research

EPSRC
Engineering and Physical Sciences
Research Council



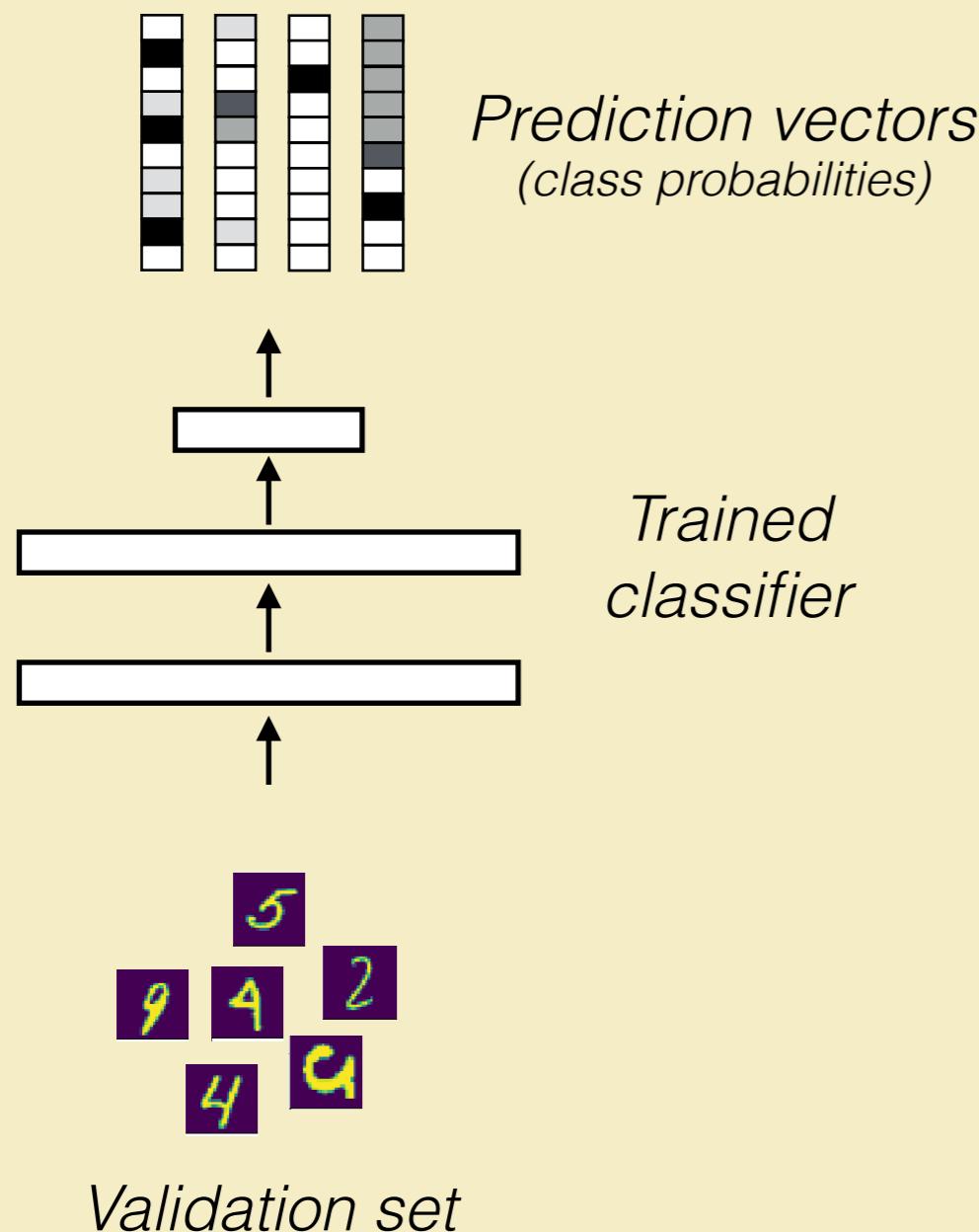
HUAWEI

BONUS! Debugging tools



Programming is to debugging
as
Differentiable programming is to
error analysis

Error Analysis via Dark Knowledge



Interpret via visualization

Visualize via dimension reduction

Dimension reduction via dark knowledge

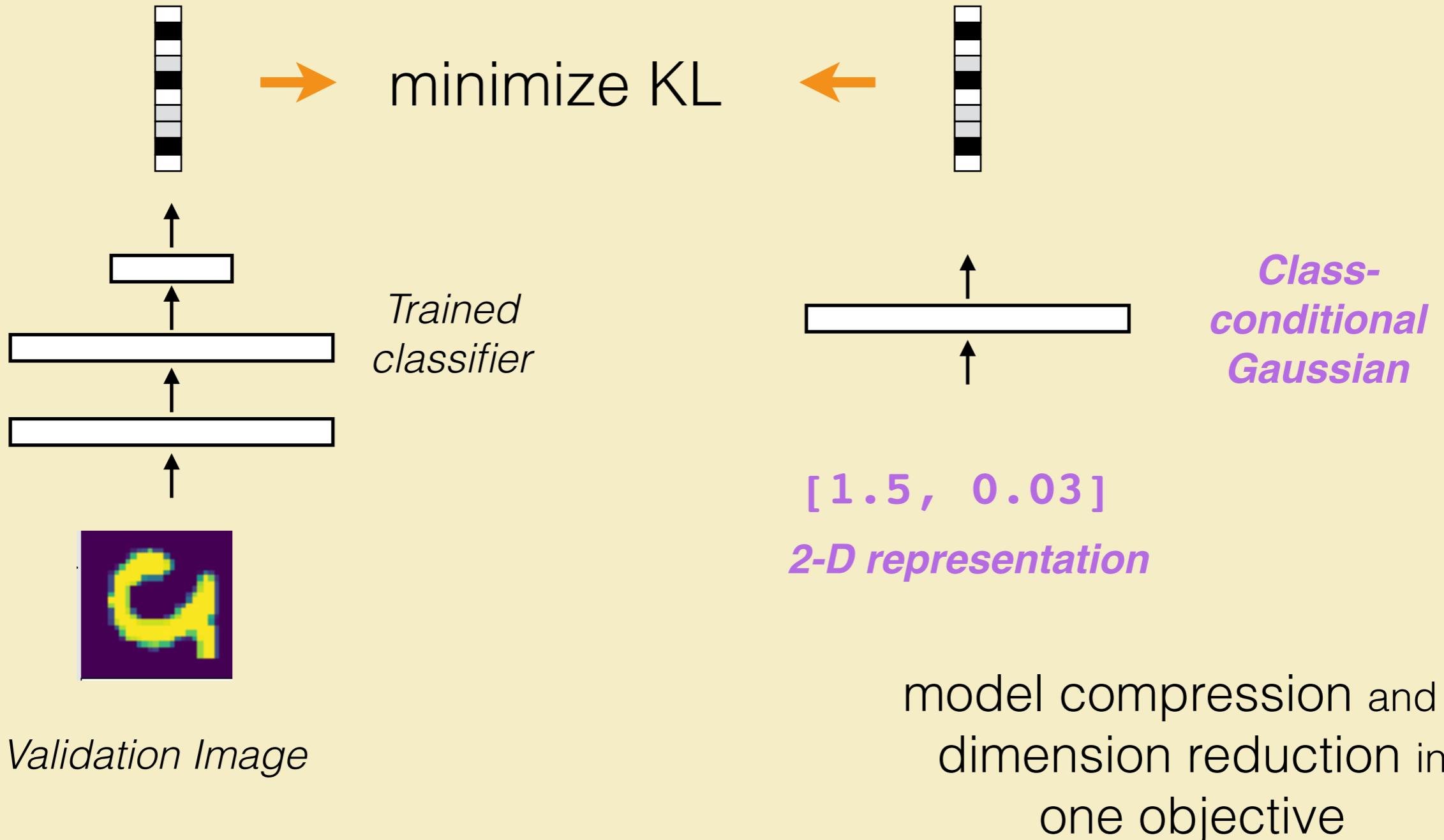
Why dark knowledge?

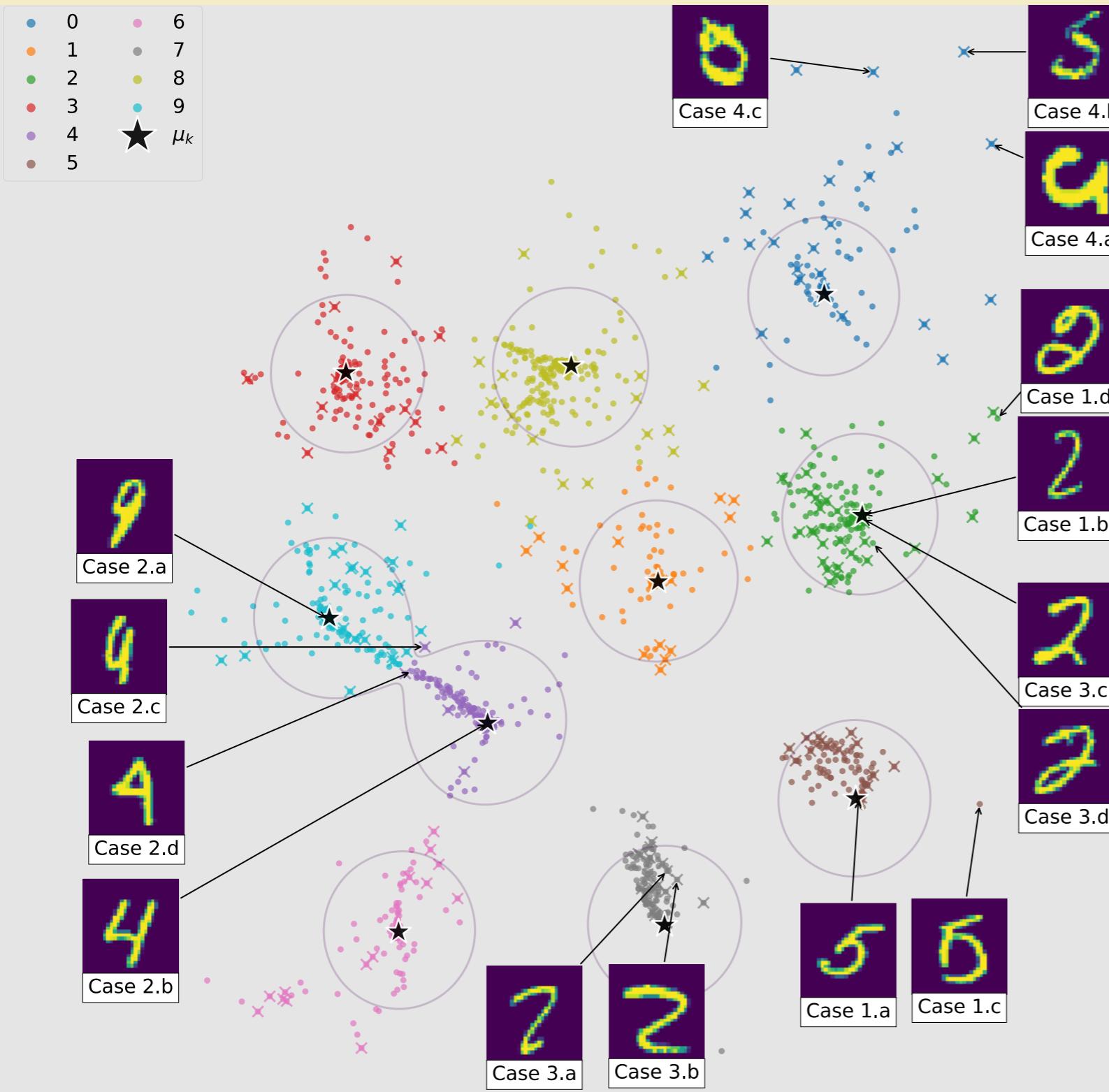
| | 0.95 | 0.95 |
|------------|------|------|
| airplane | 0.95 | 0.95 |
| bird | 0.04 | 0 |
| motorcycle | 0 | 0 |
| truck | 0 | 0 |
| frog | 0 | 0.05 |



Maybe something's wrong

Error Analysis via Dark Knowledge

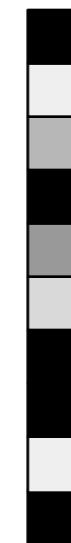




Each data item gets
2-D location
such that
Gaussian classifier 2-D
matches
original deep classifier
on original data

*Unify *this*!!!*

Mr Continuous
Representation



*Subsymbolic
son of a ...*

Professor
Symbolic AI

Great for perceptual data

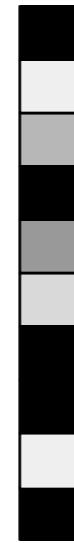
Efficient reasoning

Abstracts over lots of
perceptual states

Logical and algorithmic
reasoning

Why combine?

Mr Continuous
Representation



Professor
Symbolic AI

Great for perceptual data

Efficient reasoning

✗ Only doing local search

✗ Requires lots of data for learning

Abstracts over lots of perceptual states

Search has “non-local” effects

small change makes big difference

useful for transfer learning

✗ Impedance mismatch with perception (percepts → symbol problem)

✗ Search intractable

Why not to combine?



Over 2000 years of
overpromising and
underdelivering

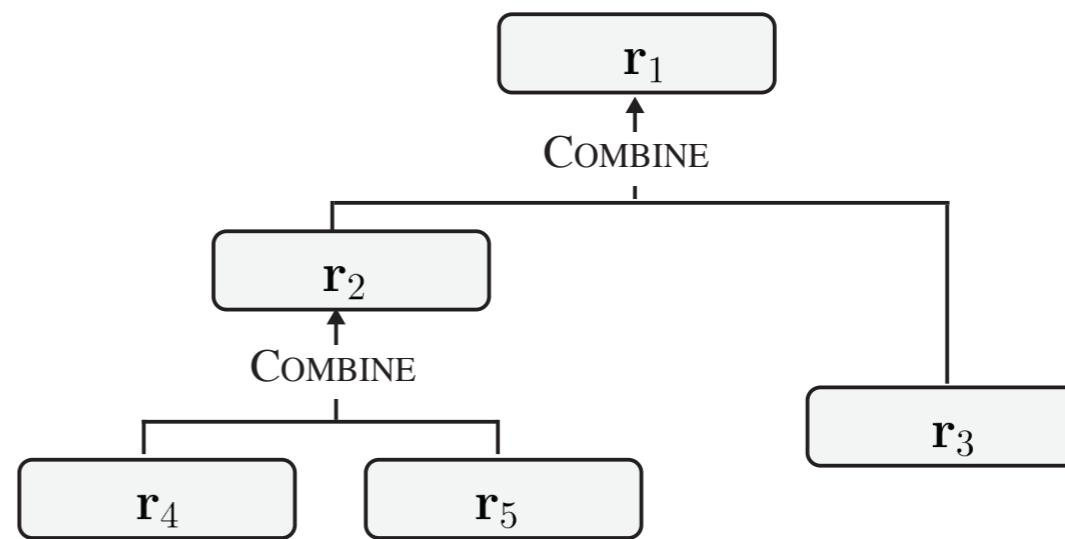
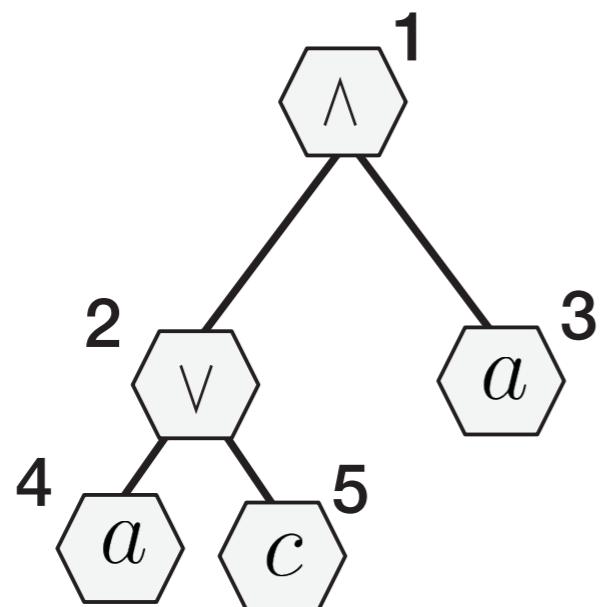
Aristotle

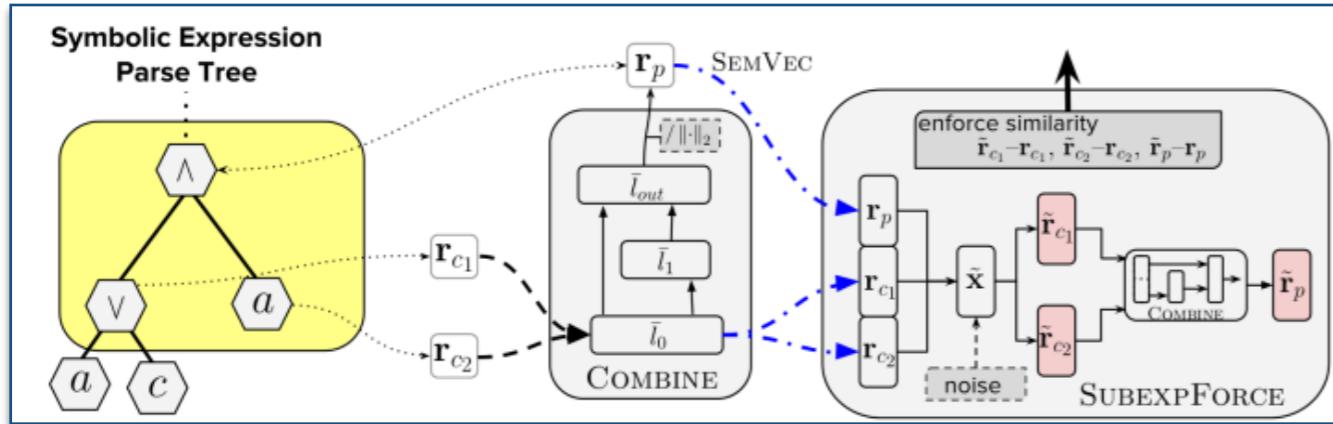
How to combine

Symbols describe structure

Use gradient to learn parameters

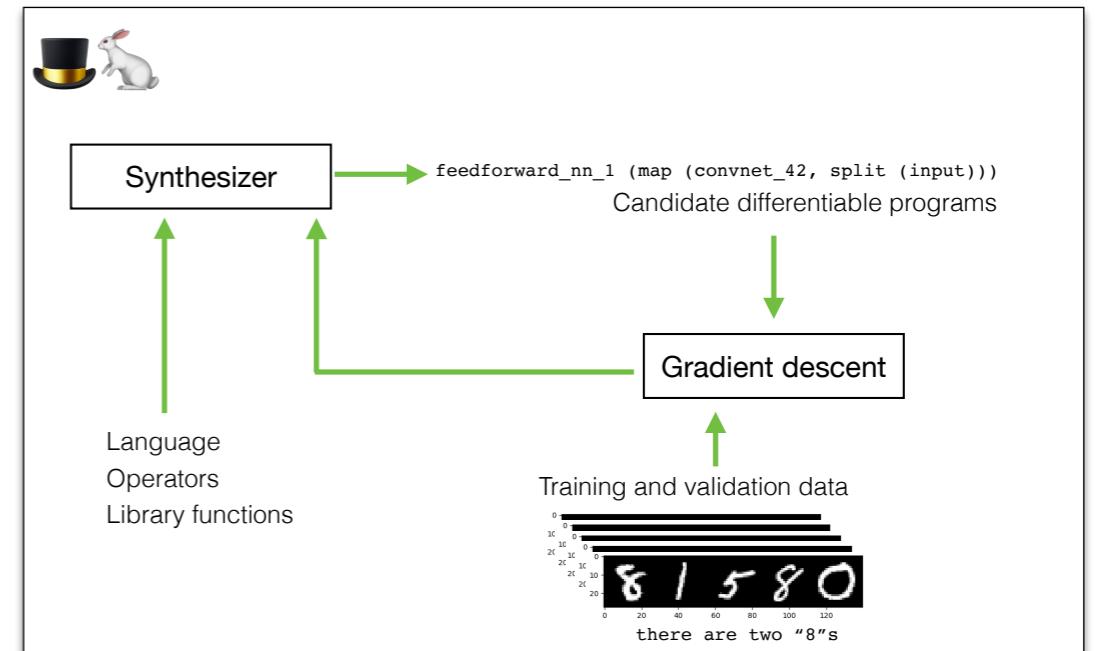
$$(a \vee c) \wedge a$$



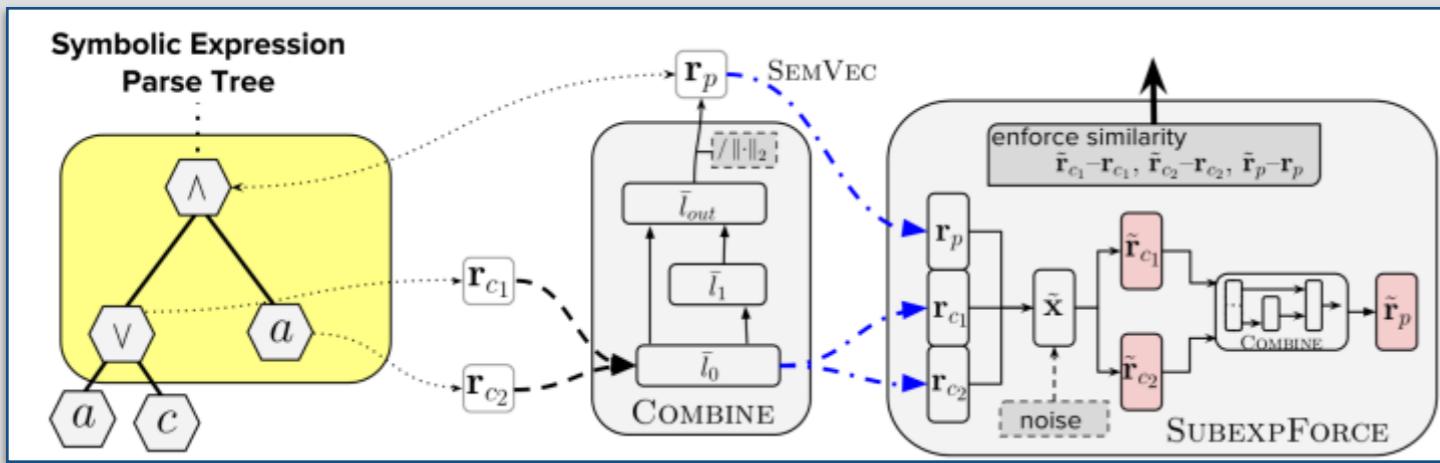


SemVecs

Representing symbolic expressions by vectors

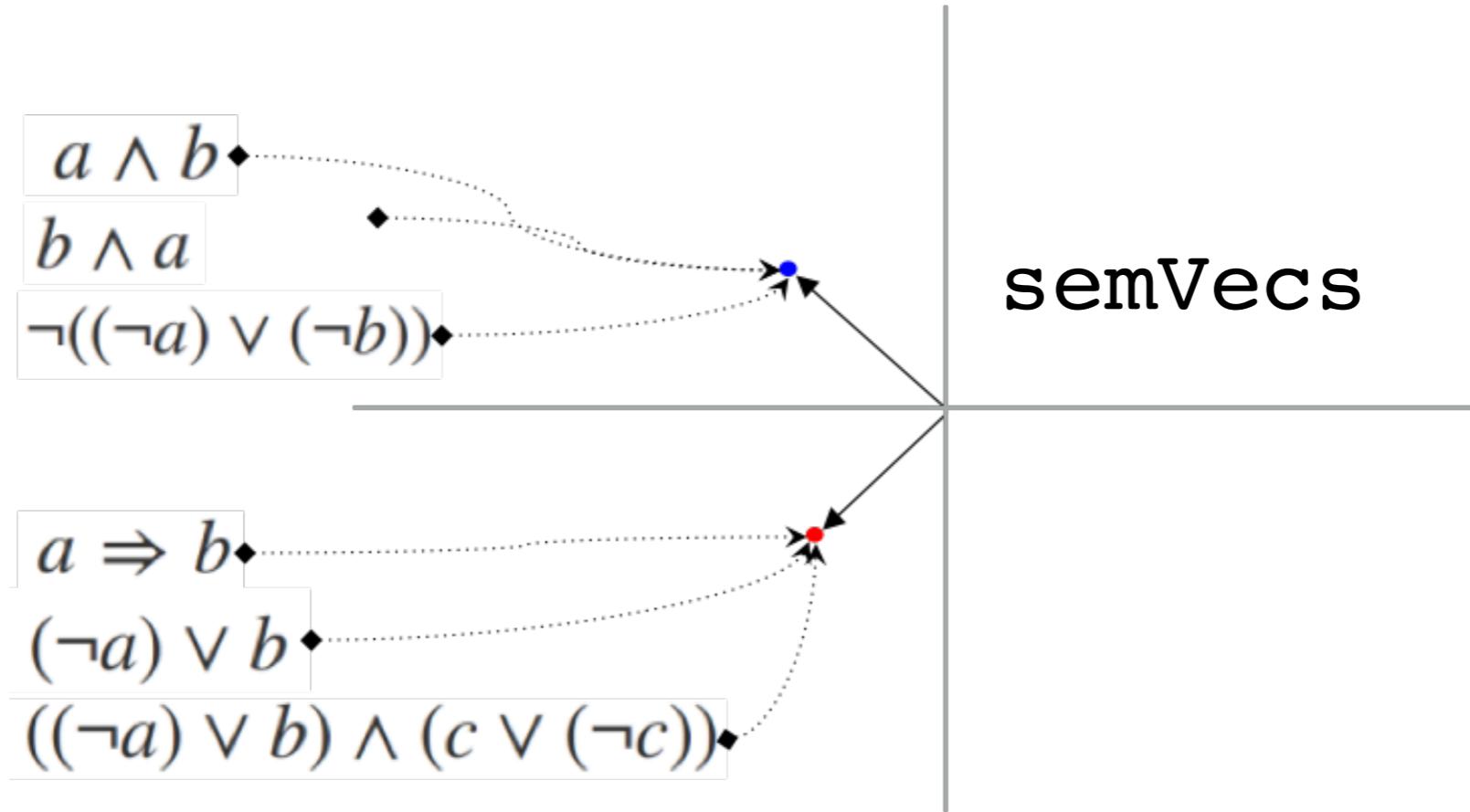


Representing differentiable functions by functional combinators



Continuous Representations of Symbolic Expressions

Can vectors help symbols?



How this works: symbol is mapped to a vector (semantic equivalence)

can we compress into **continuous** vector?

Want similar continuous vectors —> logically equivalent

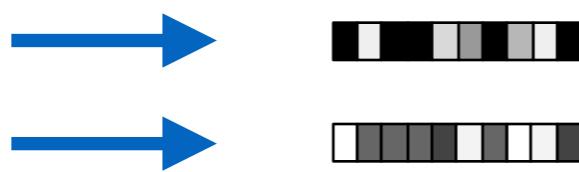
Potential Uses

Logical expressions

$$a \vee (b \implies c)$$

$$a \vee \neg b \vee c$$

Continuous vectors (`semVecs`)



Symbolic reasoning:

~~search~~

pattern recognition

Theorem Proving

[DeepMath: Irving et al, 2016]

[Zaremba et al, 2014]

Program Synthesis

[Gulwani et al, CACM 2015]

Inductive Logic Programming

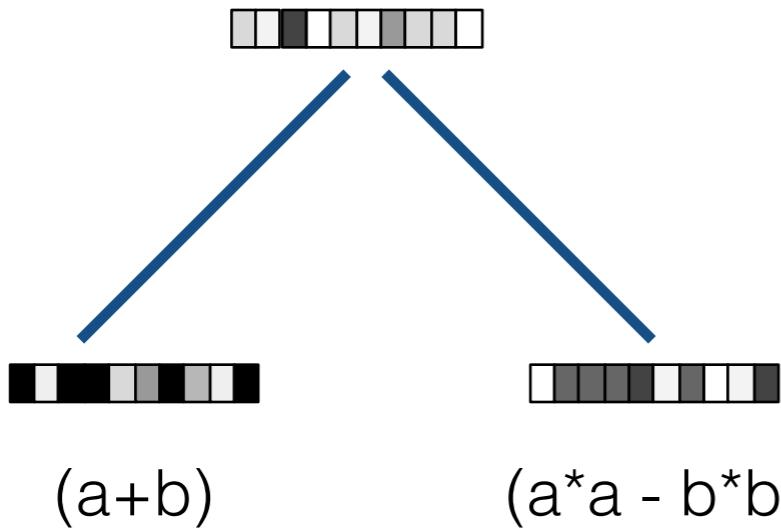
[Rocktaschel and Riedel, 2016]

[Rocktaschel and Riedel, arXiv 1705.11040 2017]

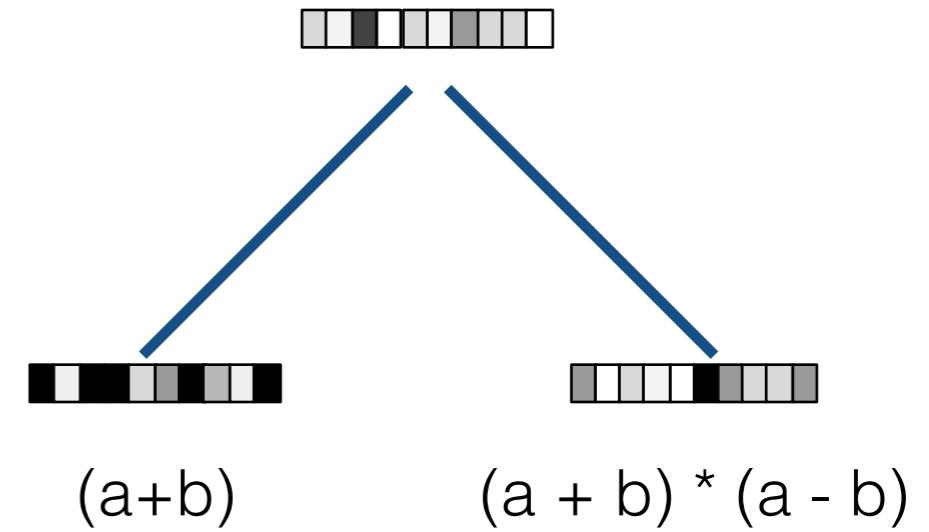
Transfer Learning

Desiderata

$$(a+b) * (a^*a - b^*b)$$



$$(a+b) * ((a+b)^* (a-b))$$

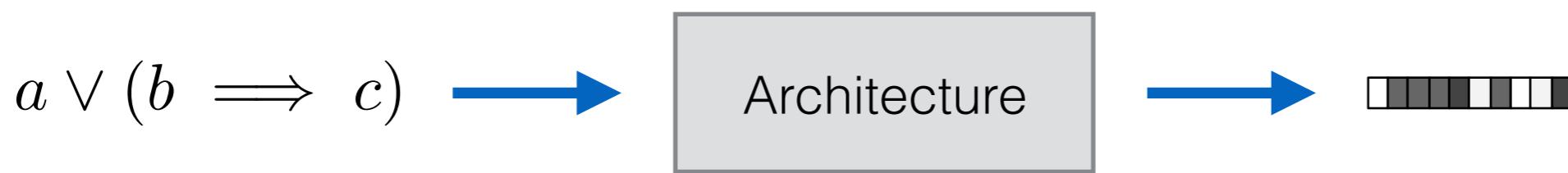


Syntax directed: Semantics is compositional

Not too much: Small syntax change → big semantics

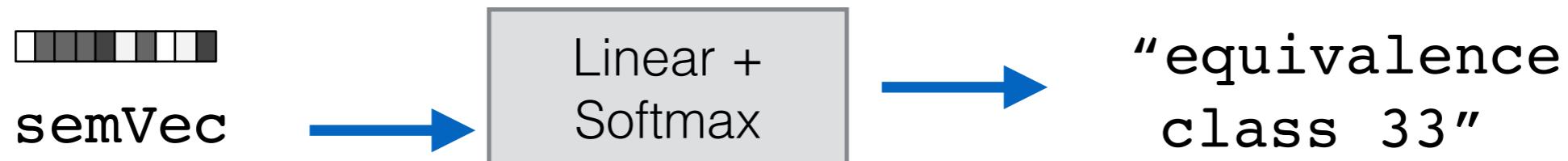
“man bites dog” problem

Computing semVecs



Training

Partition training expressions into equivalence classes



Use a supervised max-margin loss

Testing

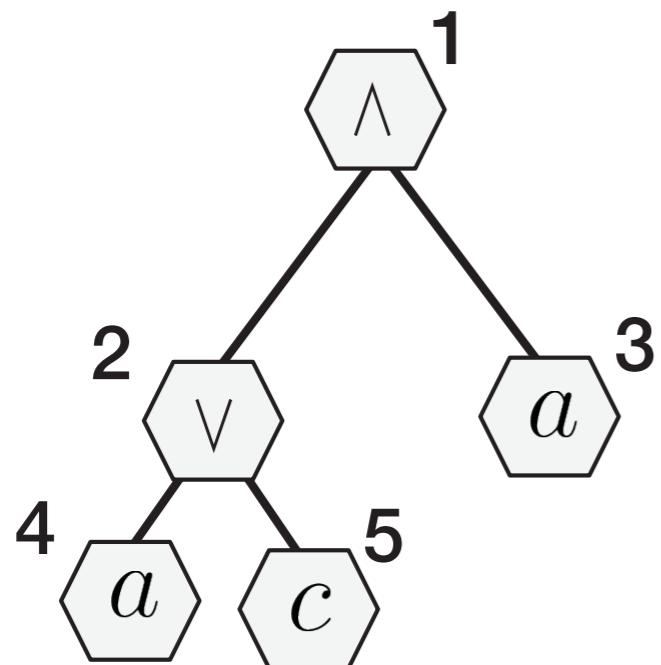
Use a semVec similarity only. Allows zero-shot learning on equiv classes.



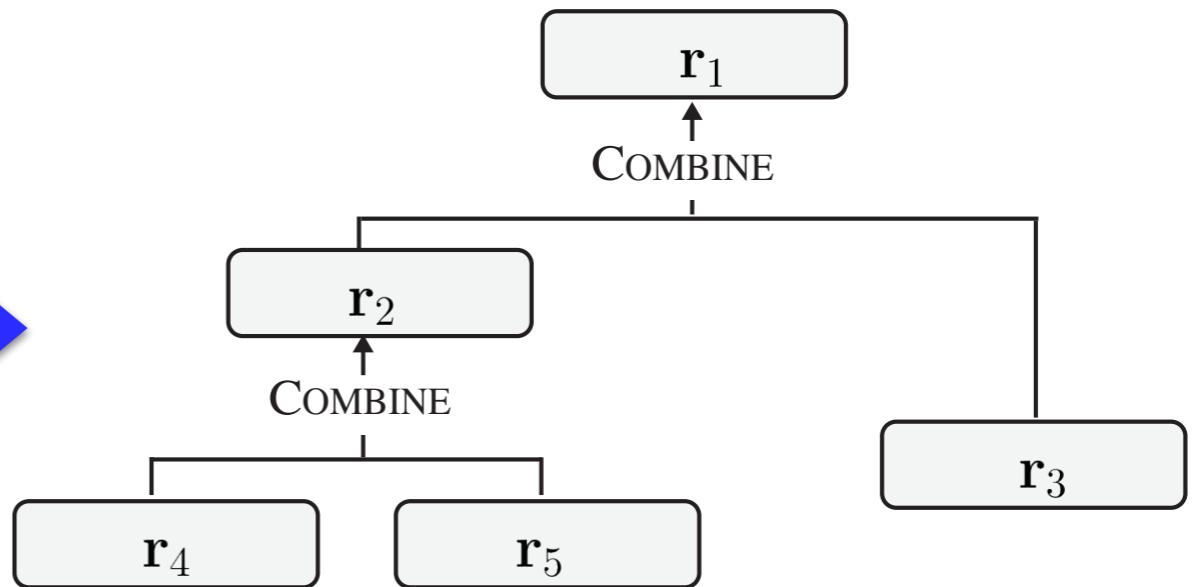
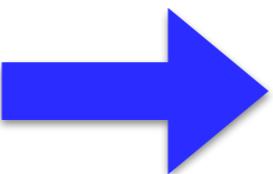
Allows zero-shot learning on equivalence classes.

Recursive NN (TreeNN)

$$(a \vee c) \wedge a$$



Syntax tree



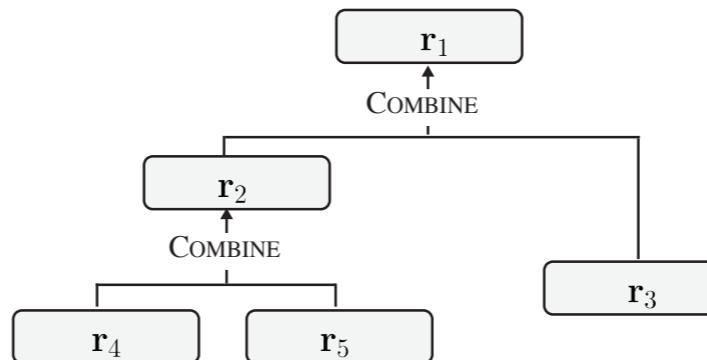
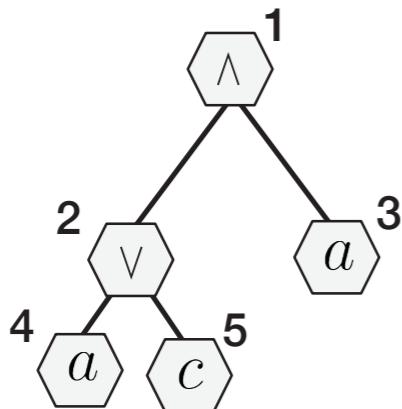
Network architecture

Problem: Representations mostly syntactic. Too much syntax!

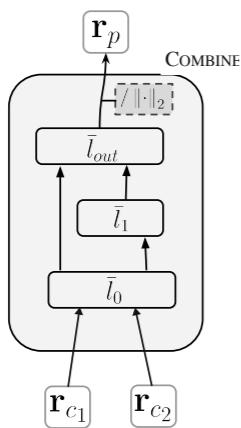
EqNet

Start with TreeNNs

$$(a \vee c) \wedge a$$



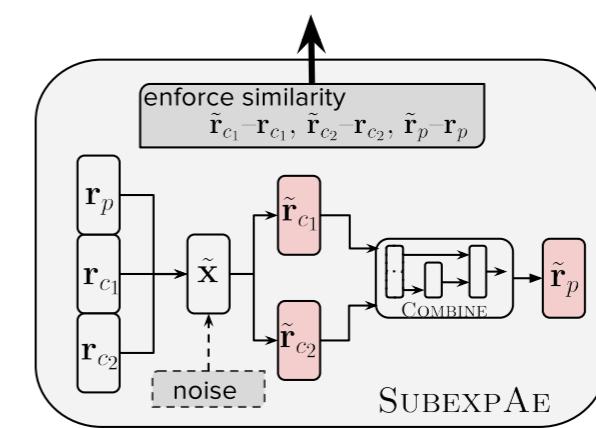
Add:



$$\|\cdot\|_2$$

Moar! Layers!

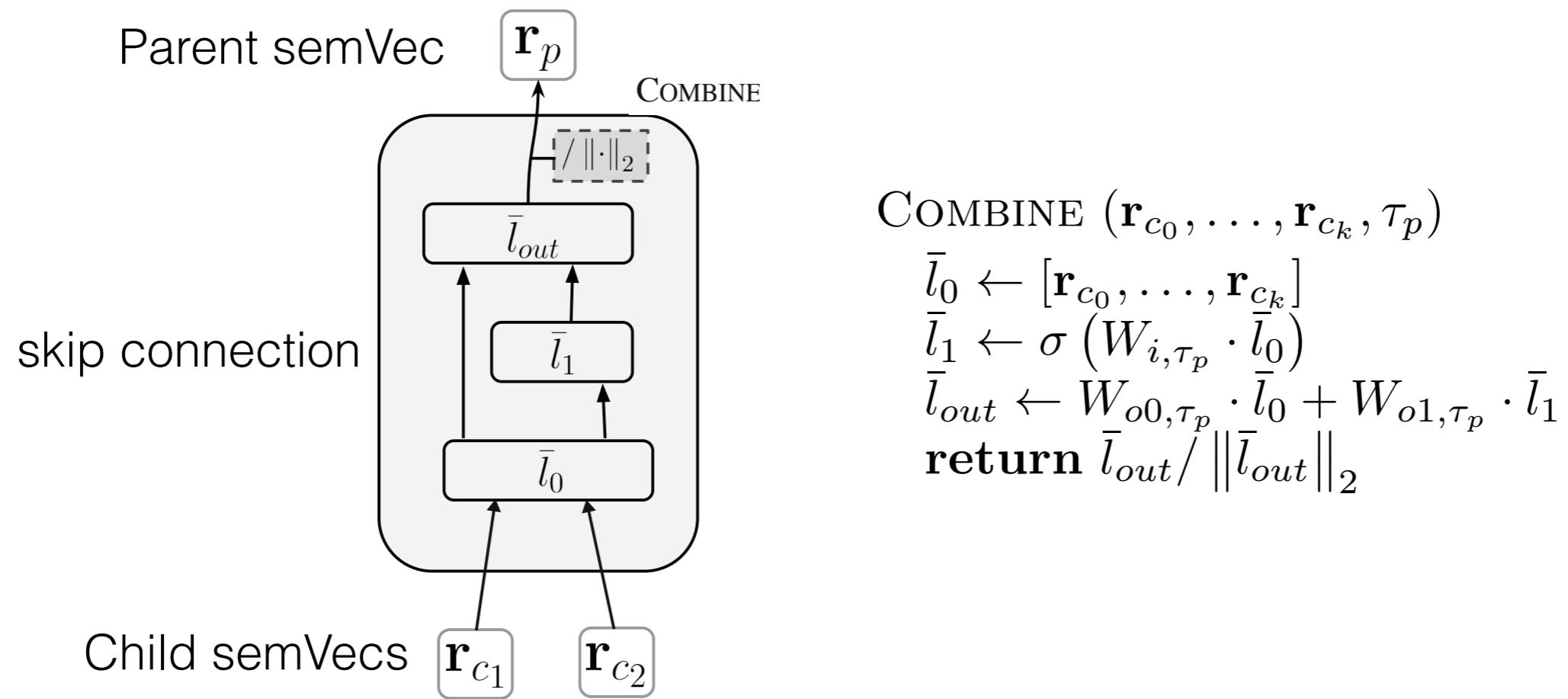
Normalization



Subexpression AE

Layers and Normalization

For one syntactic parent-child



Big impact.

(Turns out you need both residual and normalisation together)

SubexprAE: Motivation

Semantic information is bidirectional

Not only do **children** provide info re **parents**

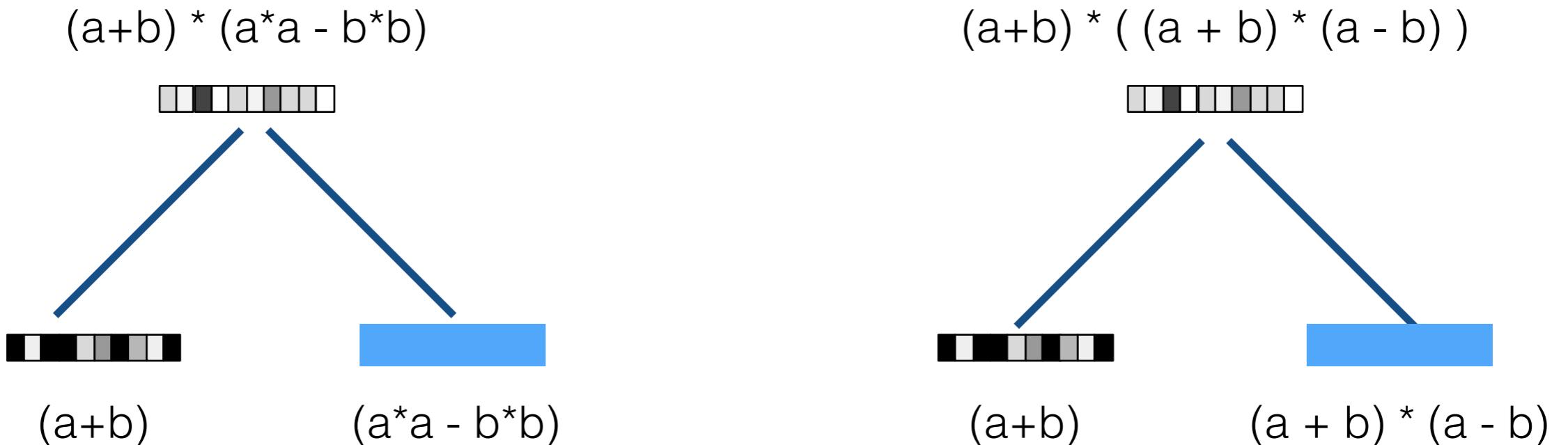
But **parents** provide info re **children**

```
uncle(?B,?A) :- parent(?Z,?A), brother(?Z,?B).
```

Unification propagates this info automatically

How to map to continuous space?

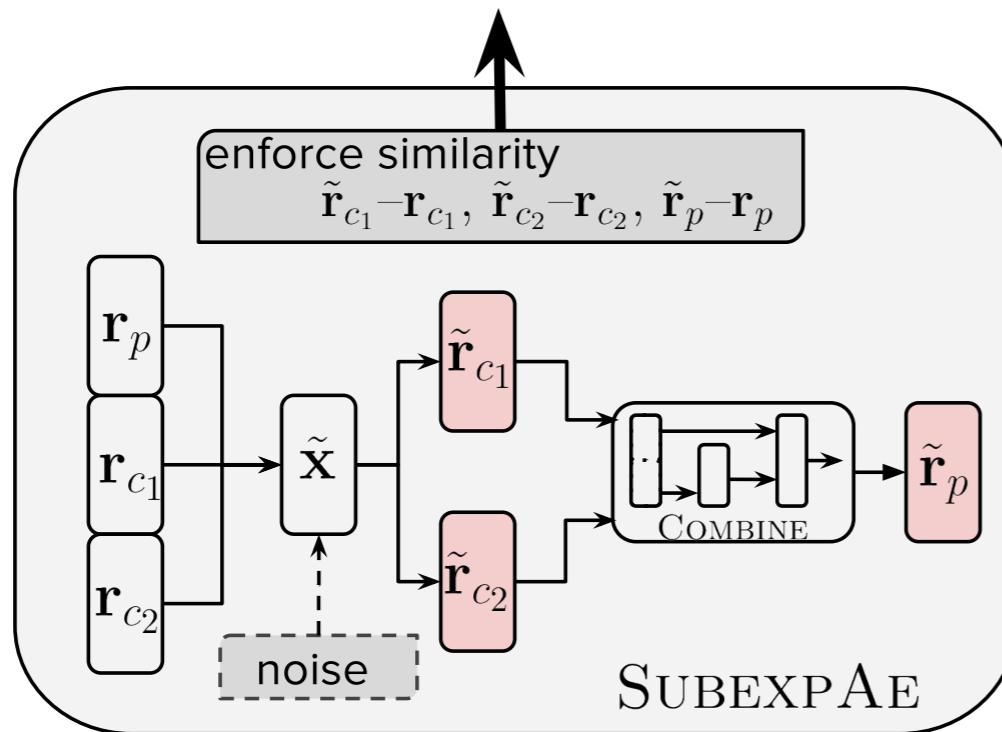
SubexprAE Motivation



ensure this prediction problem is “easy”
semantic classes will be clustered together

Subexpression Autoencoder

For every node in syntax tree, add regularisation



Denoising autoencoder
plus bottleneck on
(parent, child1, child2)
semVecs

Intention is

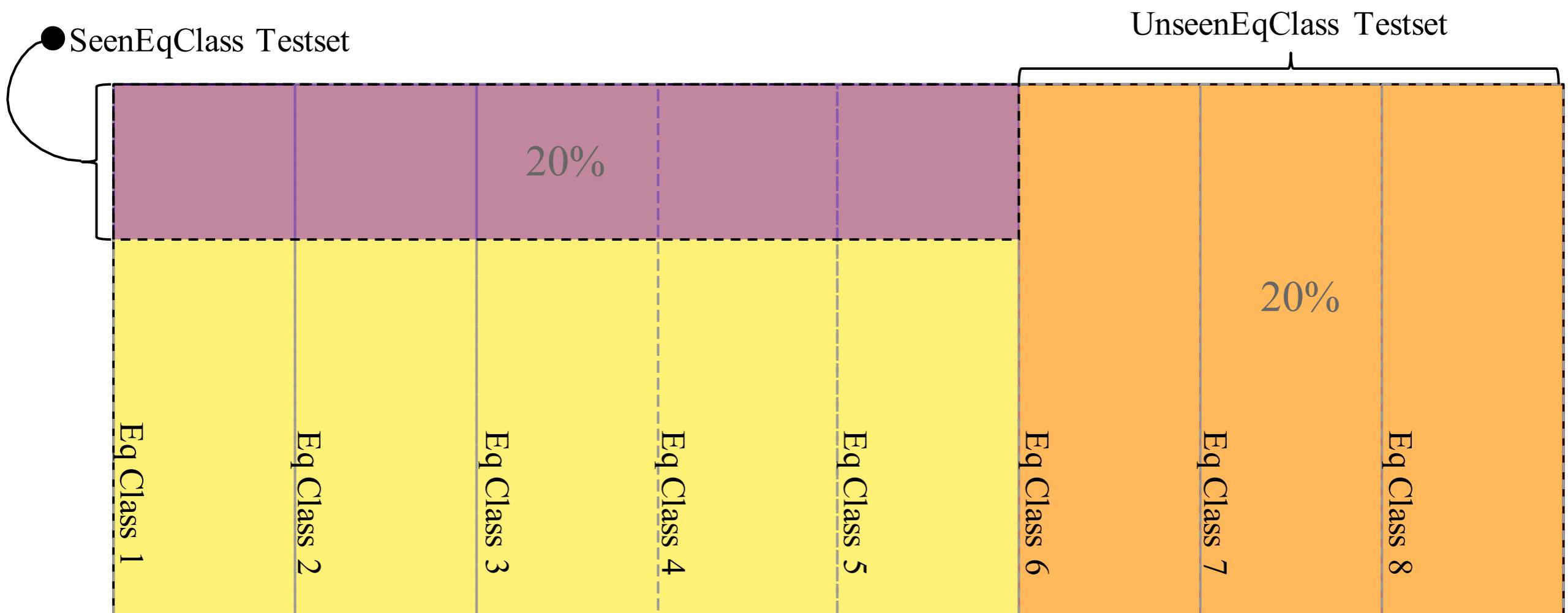
Bottleneck \rightarrow Abstraction

Denoising \rightarrow Reversibility

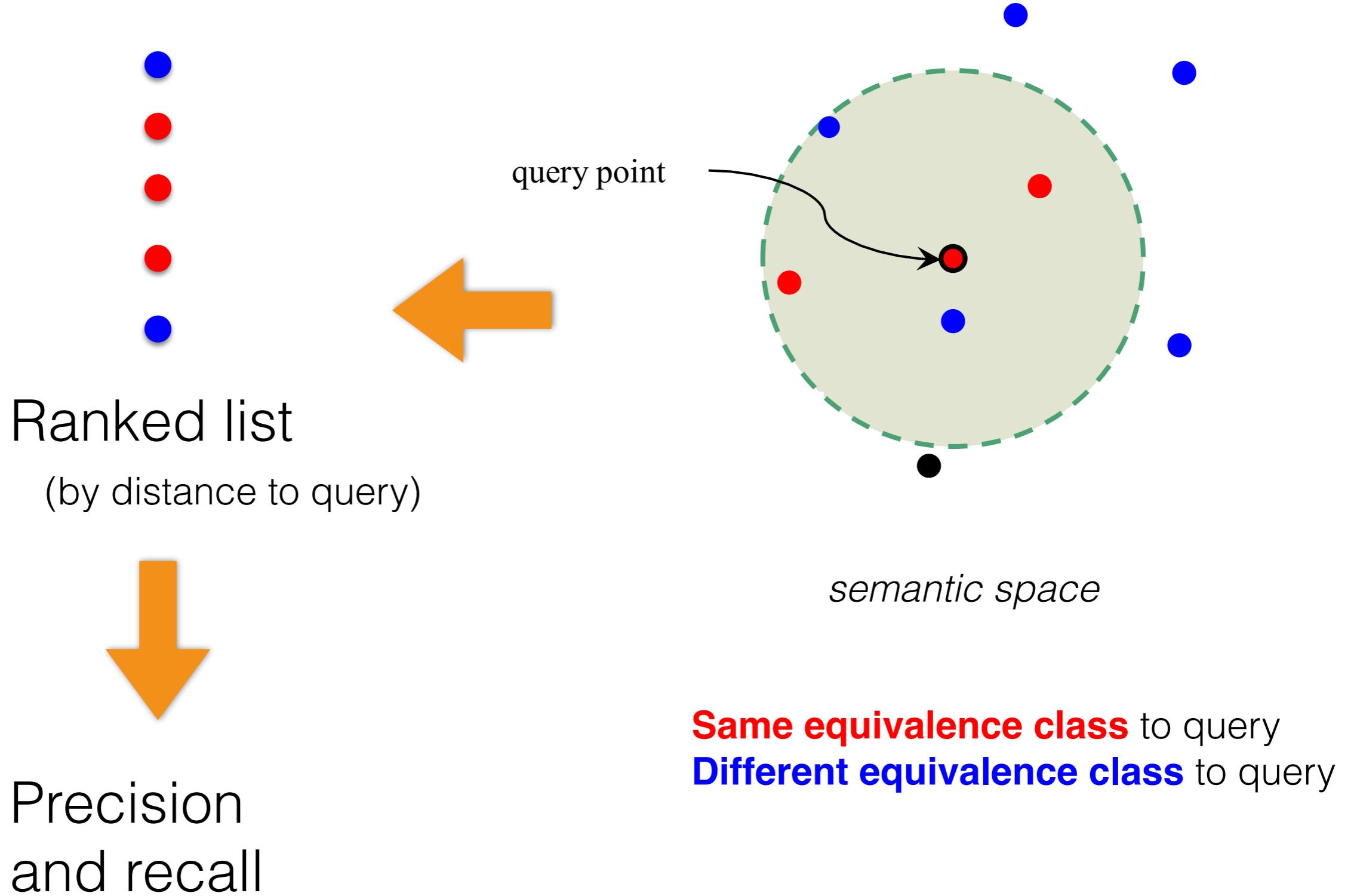
Evaluation

| Dataset | # Vars | # Equiv Classes | # Exprs | H |
|-------------------------|-----------|--------------------|------------|------|
| SIMPBOOL8 | 3 | 120 | 39,048 | 5.6 |
| SIMPBOOL10 ^S | 3 | 191 | 26,304 | 7.2 |
| BOOL5 | 3 | 95 | 1,239 | 5.6 |
| BOOL8 | 3 | 232 | 257,784 | 6.2 |
| BOOL10 ^S | 10 | 256 | 51,299 | 8.0 |
| SIMPBOOLL5 | 10 | 1,342 | 10,050 | 9.9 |
| BOOLL5 | 10 | 7,312 | 36,050 | 11.8 |
| SIMPPOLY5 | 3 | 47 | 237 | 5.0 |
| SIMPPOLY8 | 3 | 104 | 3,477 | 5.8 |
| SIMPPOLY10 | 3 | 195 | 57,909 | 6.3 |
| ONEV-POLY10 | 1 | 83 | 1,291 | 5.4 |
| ONEV-POLY13 | 1 | 677 | 107,725 | 7.1 |
| POLY5 | 3 | 150 | 516 | 6.7 |
| POLY8 | 3 | 1,102 | 11,451 | 9.0 |

Training / Test Split

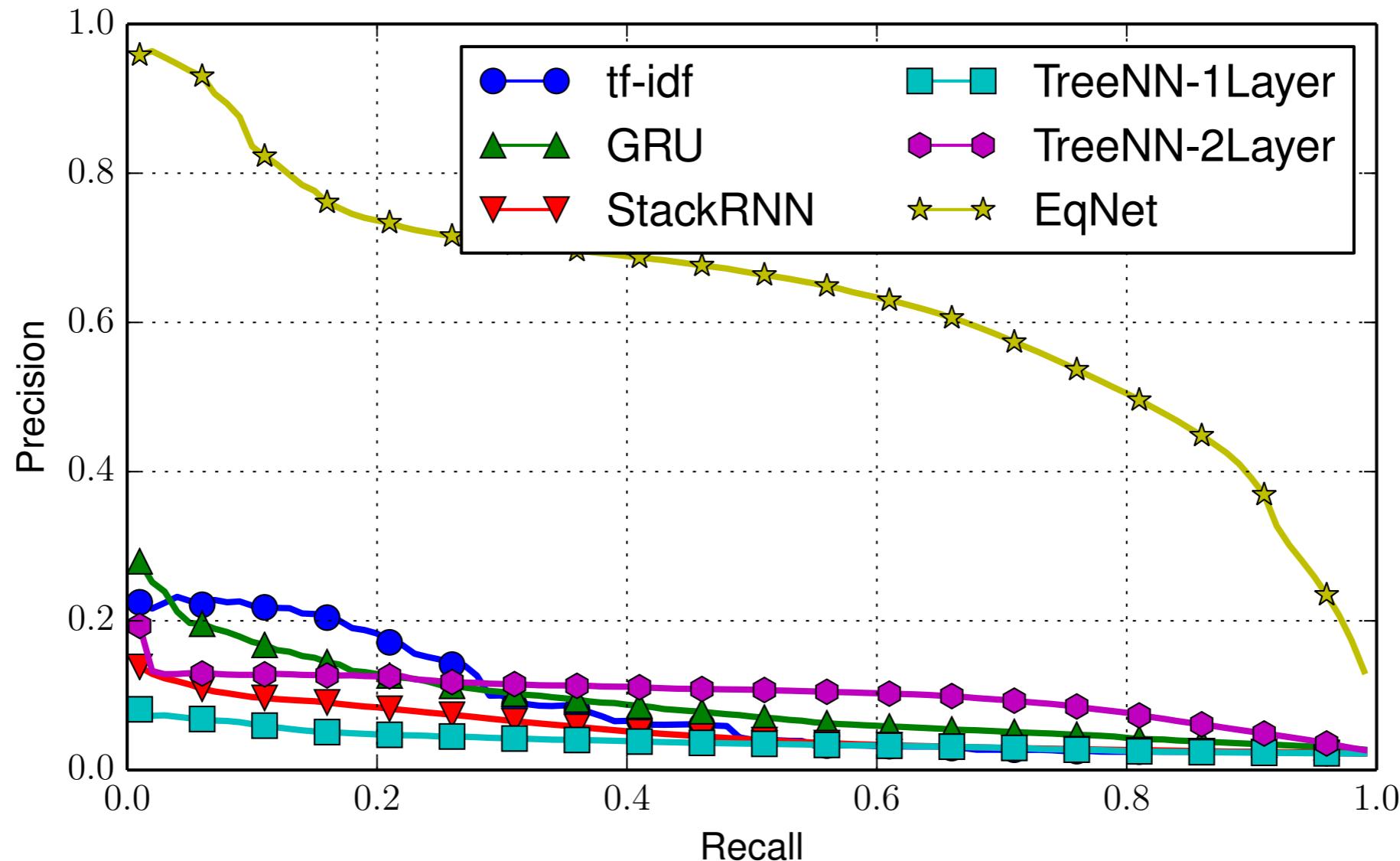


Evaluation Metric



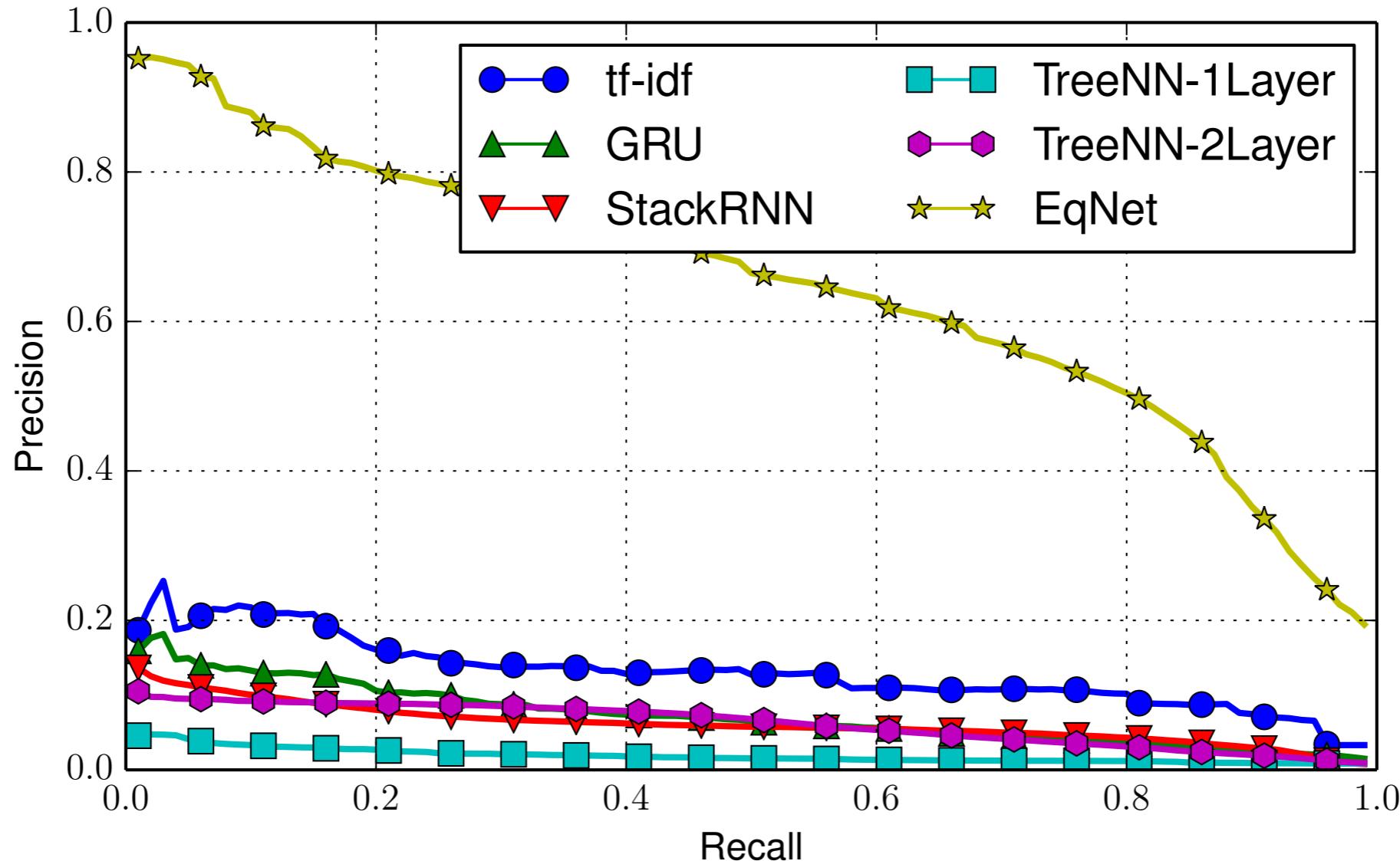
Seen equivalence classes

Equivalent expressions to the queries were in training set



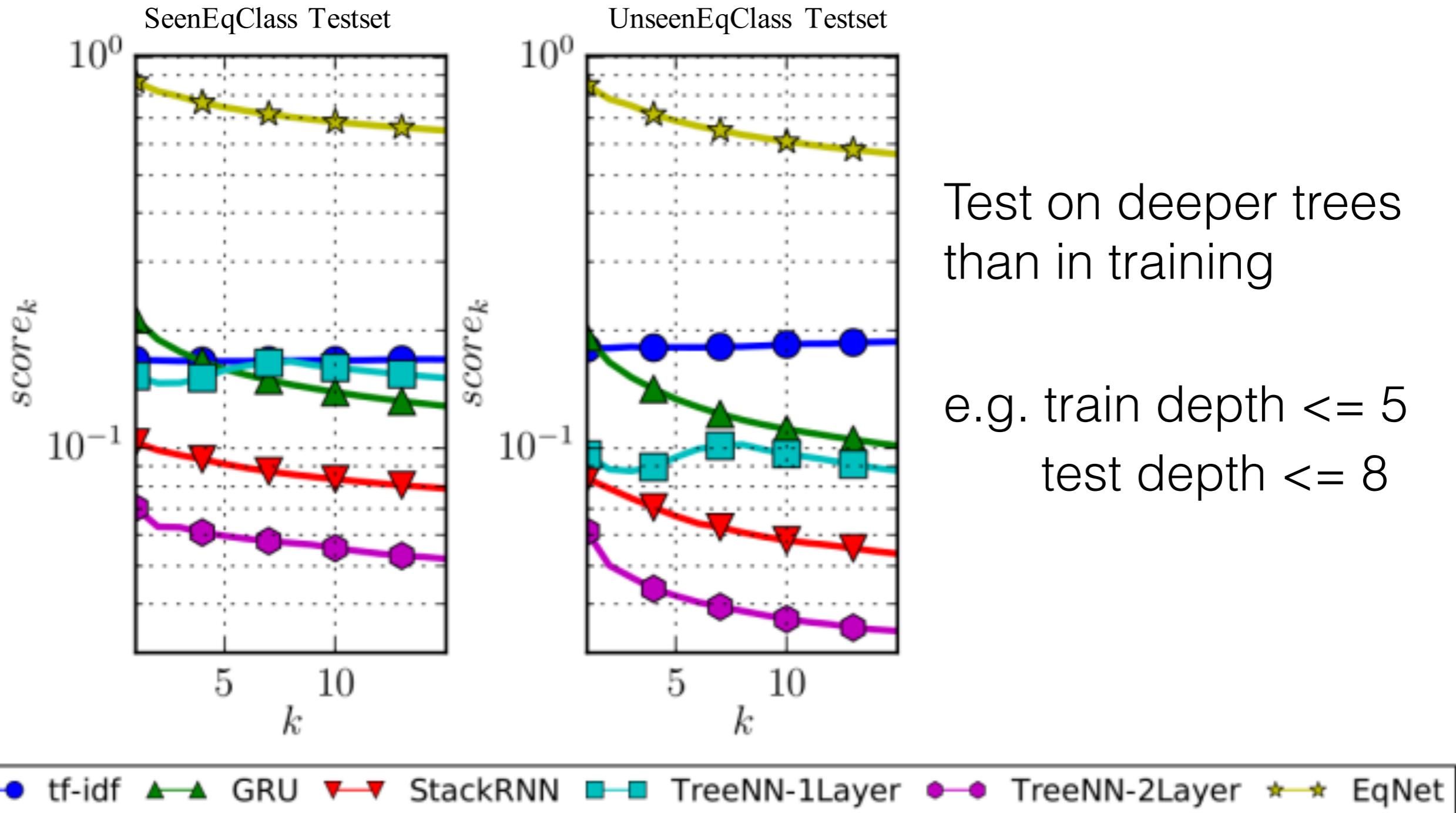
Unseen equivalence classes

Zero shot learning. No training examples of equivalent expressions.



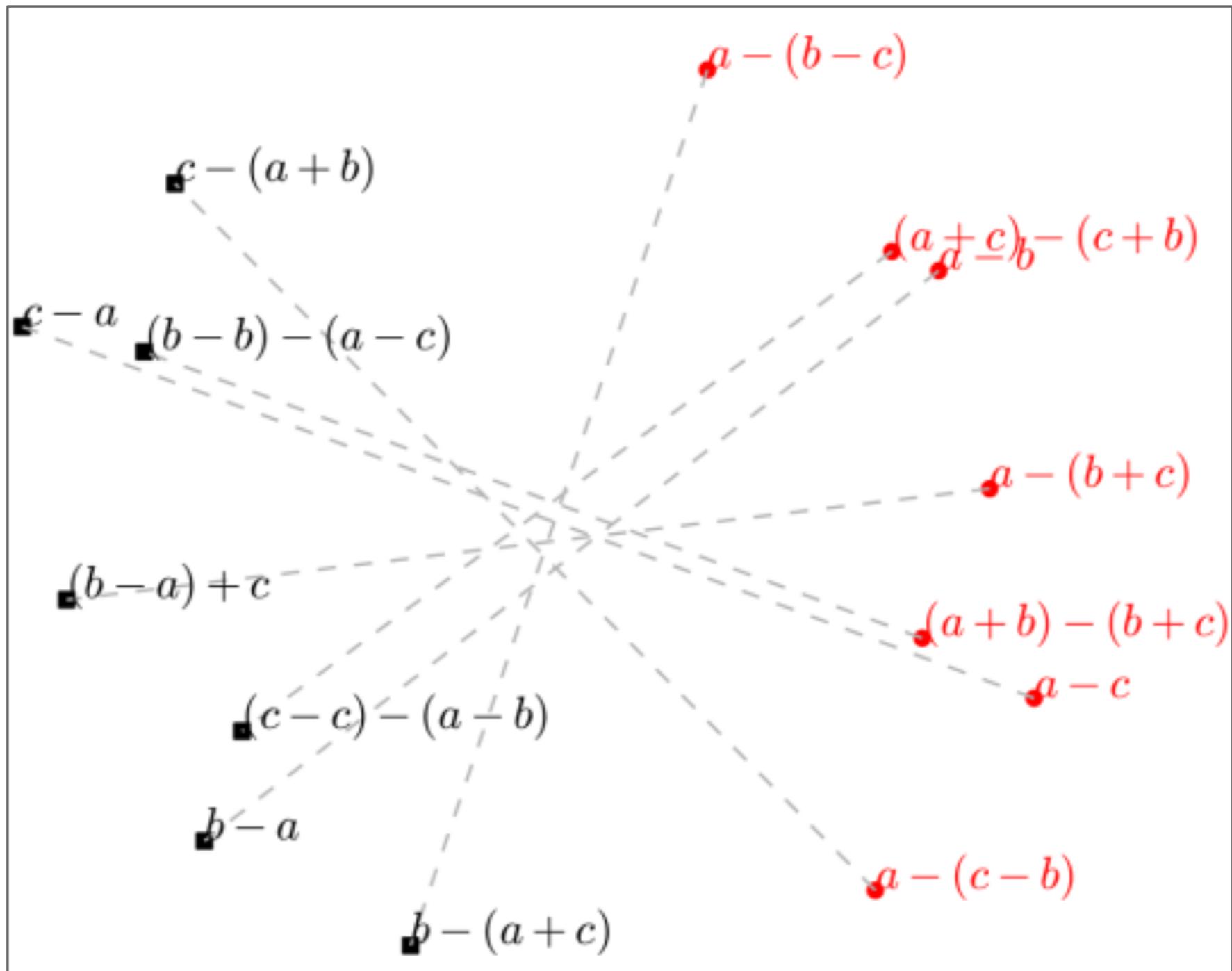
EqNet performance on seen and unseen is similar!

Learned compositionality?



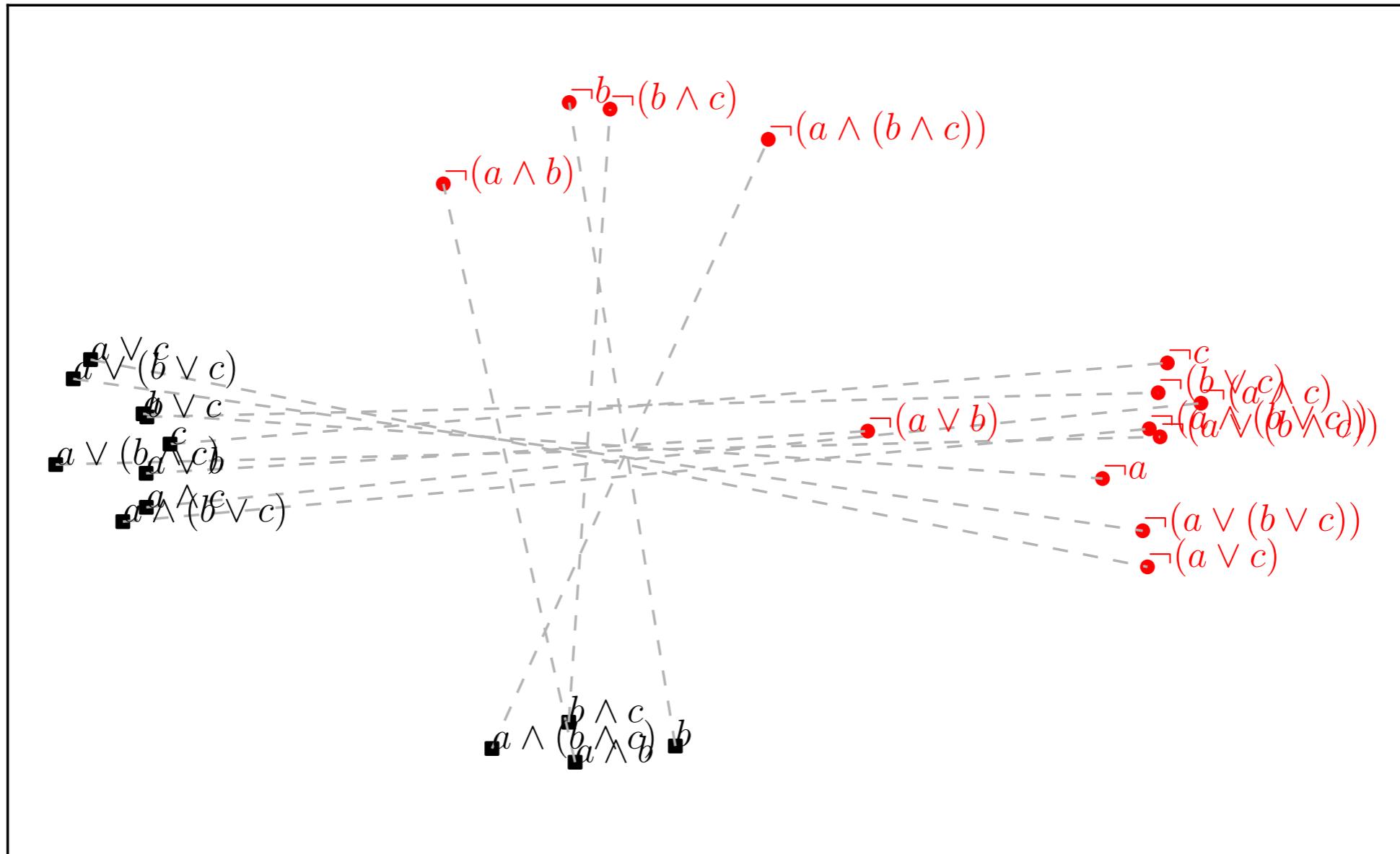
Visualizing polynomials

multivariatePolynomial2vec?

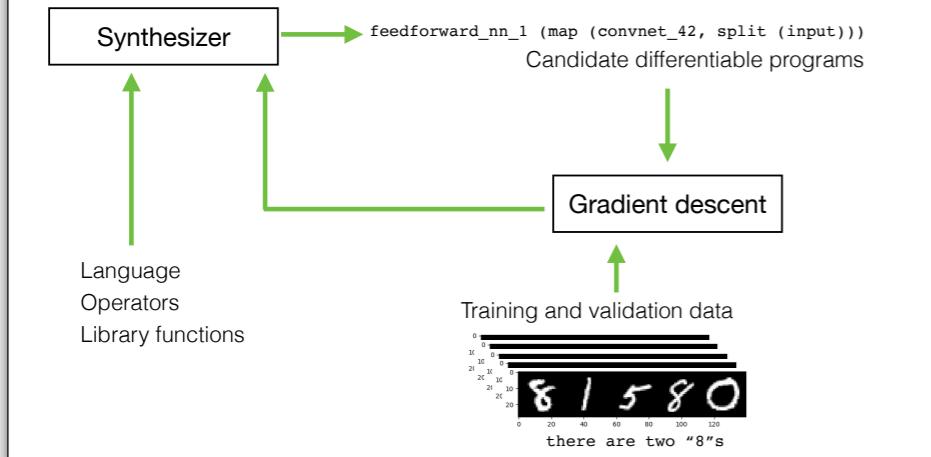


Visualizing boolean expression

booleanExpression2vec?



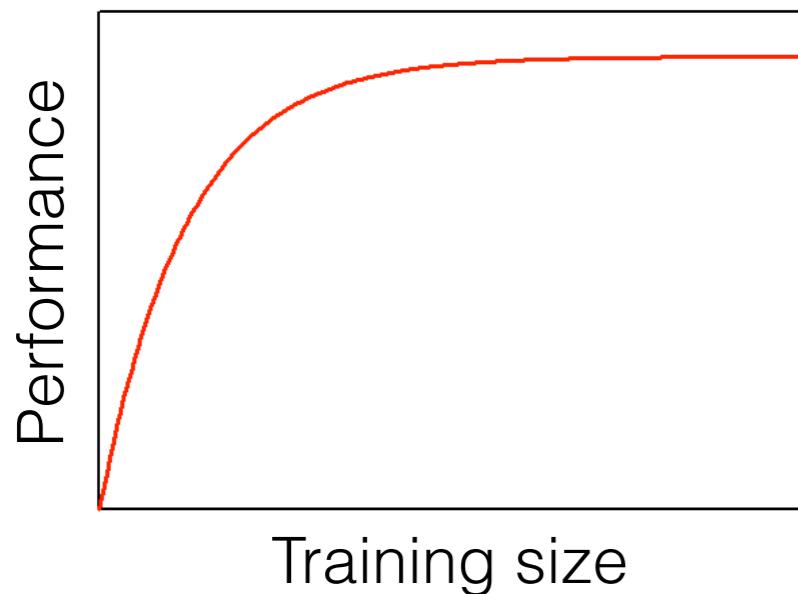
PCA visualization of semVecs



Synthesis of Differentiable Functional Programs for Lifelong Learning

[Valkov, Chaudhari, Srivastava, Sutton, and Chaudhuri,
bit.ly/adventures-neurosymbolic 2018]

The Problem with Learning



All learning curves stop.

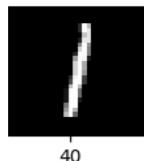
Why “data hungeriness”?

Learning systems never get better at learning itself

Lifelong learning [*Thrun & Mitchell, 1995; Carlson et al, 2010*]

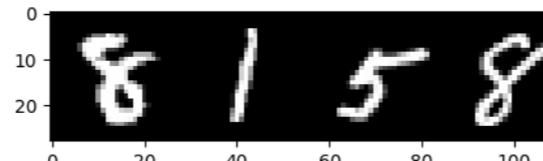
On series of tasks, accelerate learning by re-using learned structure

Task 1
Recognition



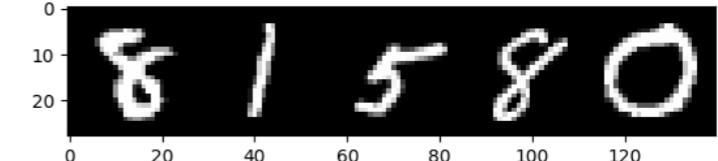
that's a “one”

Task 2
Counting



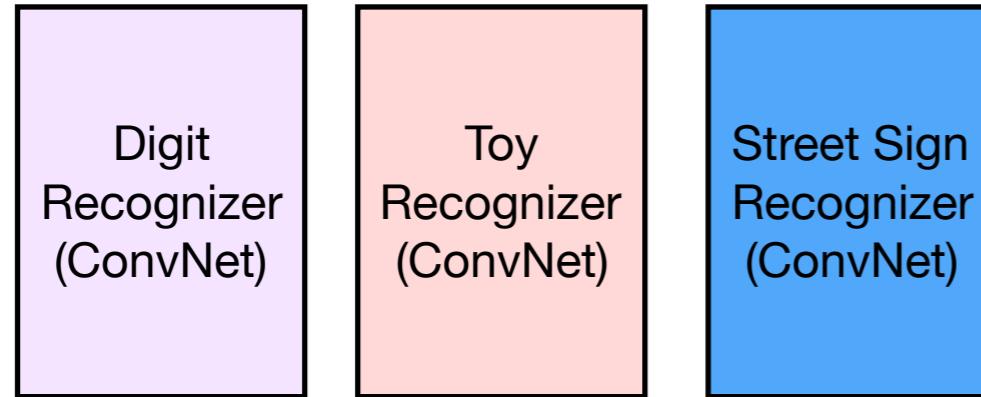
I count two “8’s”

Task 3
Arithmetic



The sum is 23

Neural libraries



Networks from old tasks

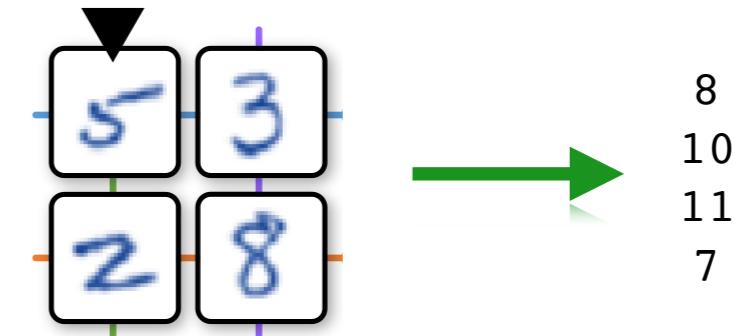
Controller in
differentiable PL



Differentiable
interpreter



New task



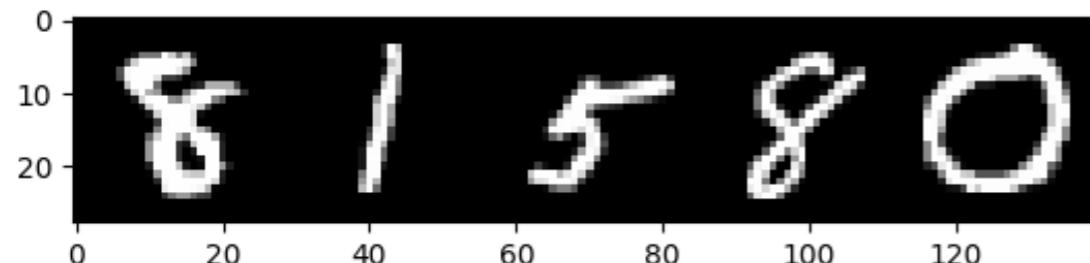
Learn programs plus
perceptual networks

```
# initialization:  
R0 = READ  
# program:  
R1 = MOVE_EAST  
R2 = MOVE_SOUTH  
R3 = SUM(R0, R1)  
R4 = NOOP  
return R3
```

end-to-end
gradient descent

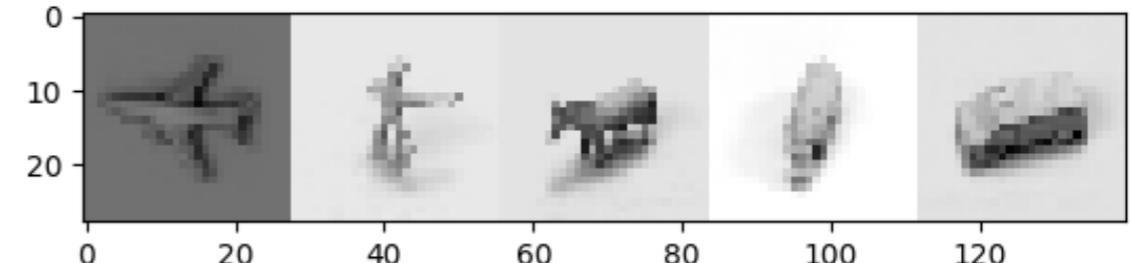
High level transfer

Task 1
Count digits

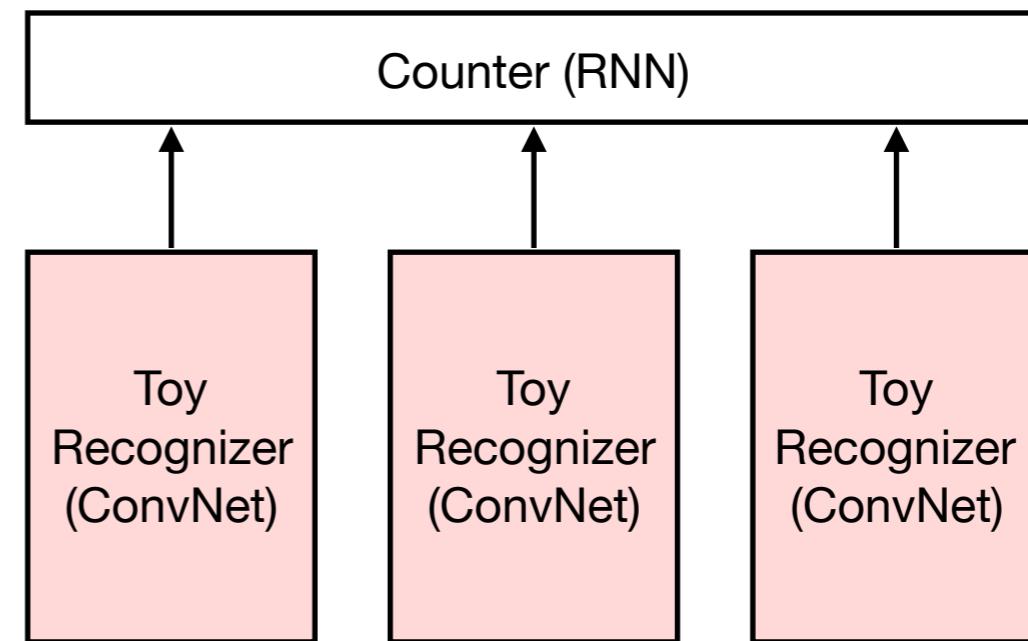
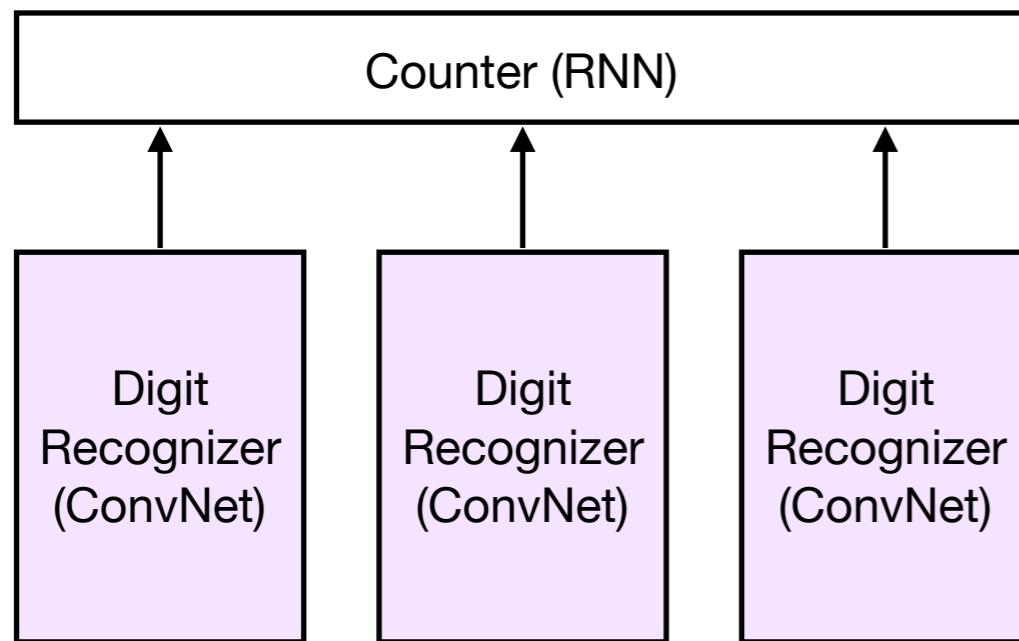


there are two "8"s

Task 2
Count toys



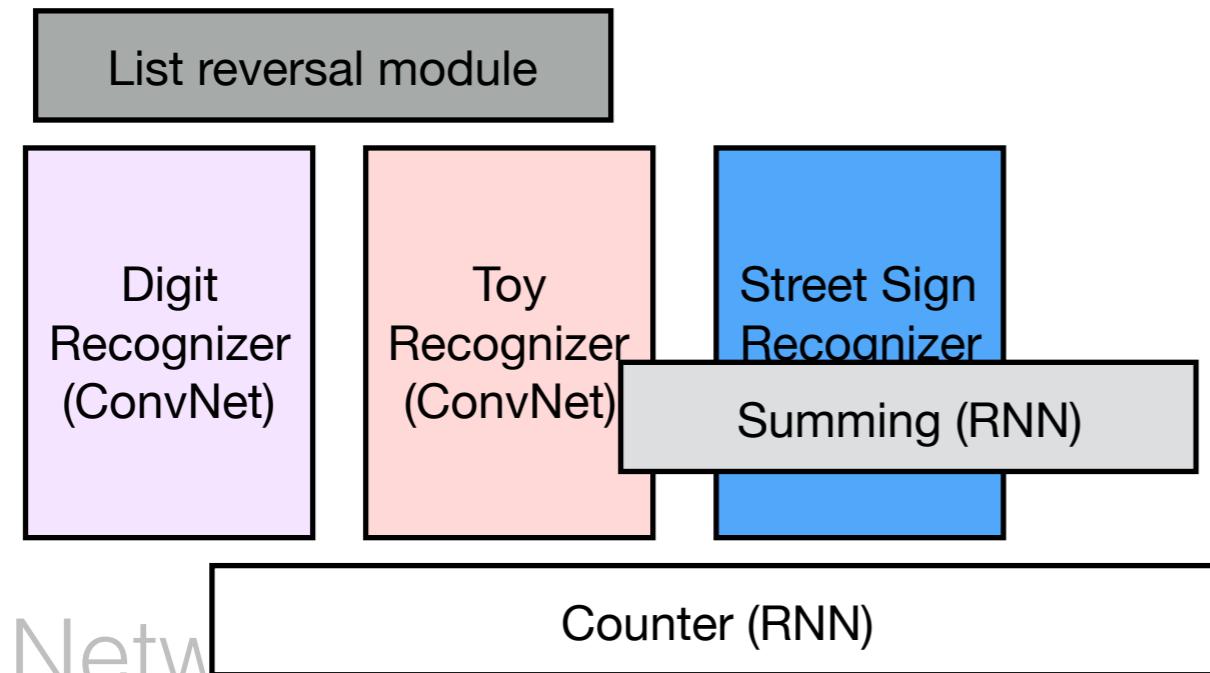
there is one "toy airplane"
(and why don't I have two?)



Reusing early layers not sufficient!

[Hinton & Salakhutdinov, 2006; Rusu et al 2016]

High-level neural libraries



Controller in
differentiable PL



Differentiable
interpreter



Problem: How to combine?

Learn programs plus
perceptual networks

```
# initialization:  
R0 = READ  
# program:  
R1 = MOVE_EAST  
R2 = MOVE_SOUTH  
R3 = SUM(R0, R1)  
R4 = NOOP  
return R3
```

end-to-end
gradient descent

New task

"High-level modules"
Other networks as input

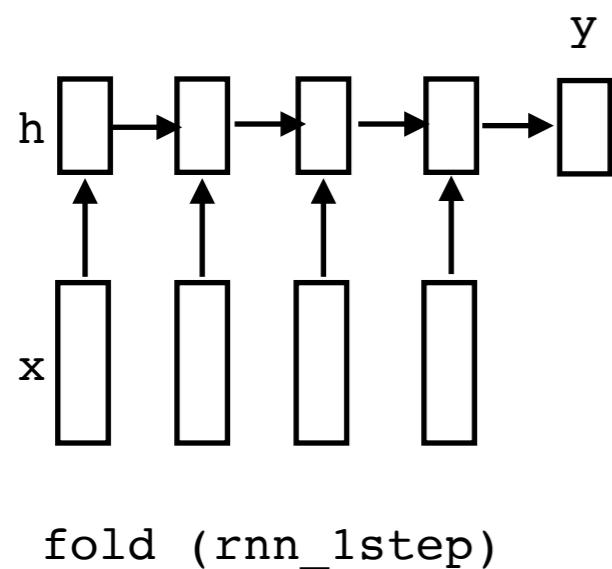


11
7

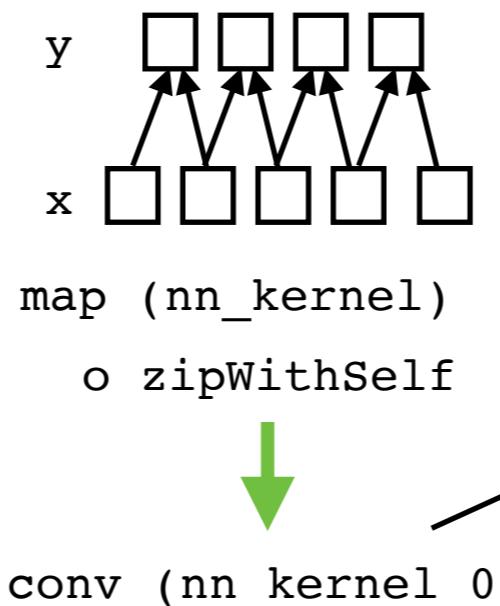
100 neural architectures, 1 weird trick

Functional programming

Recurrent neural network



1-layer ConvNet



*Combinators preserve
differentiability*

*Often point-free
[Backus, 1978]*

multiple filters? change this

Deep feedforward net

`relu o w_n o relu o ... o relu o w_1`

Deep ConvNet

`conv (nn_kernel_0) o conv (nn_kernel_1) o ... o conv (nn_kernel_D)`

Attention mechanism

`softmax o map(attn_network) o fold(rnn_1step)`

Graph convolutions

`gconv (nn_kernel_0)`

Main ideas

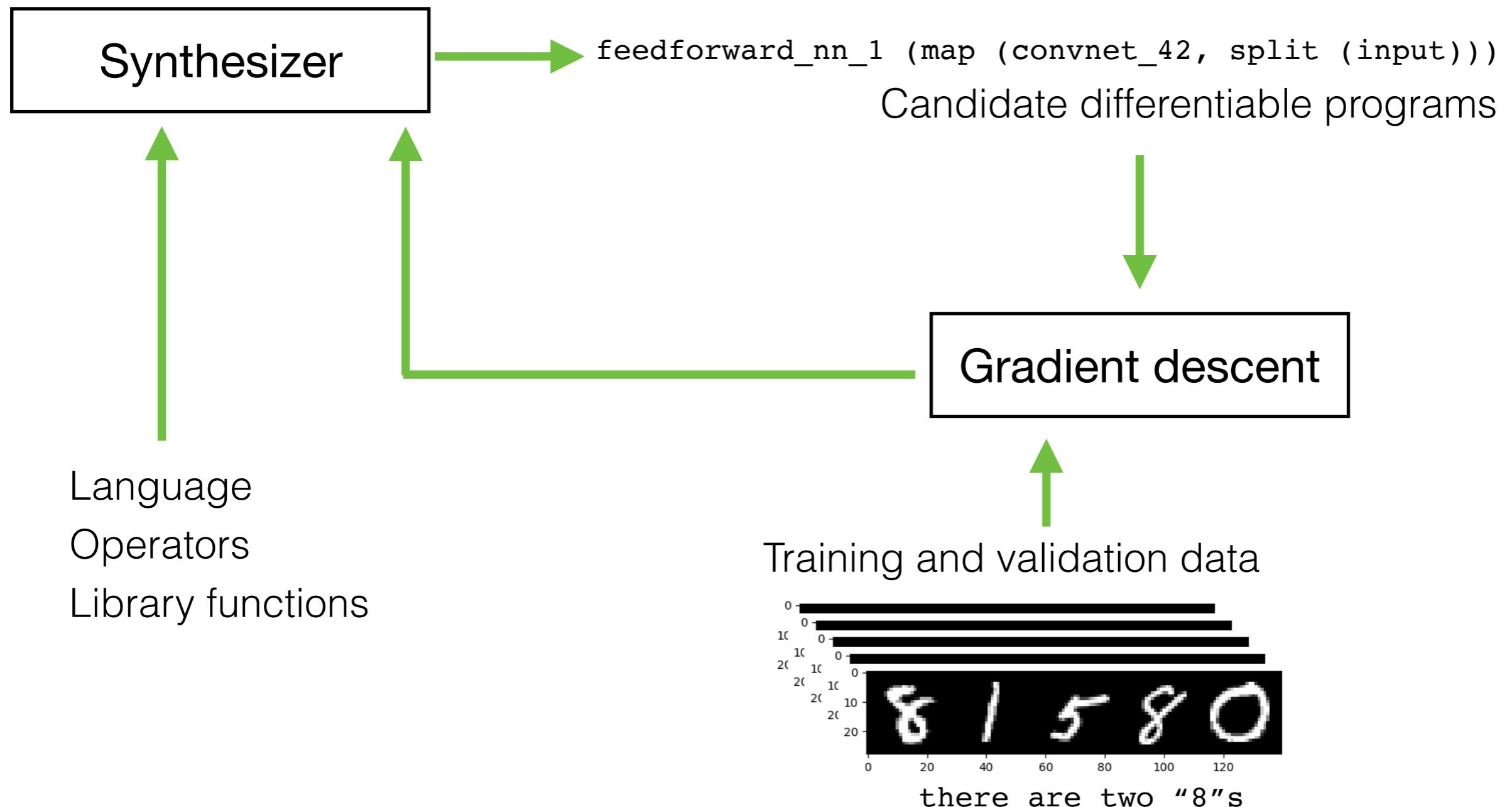
- Neural libraries for high level transfer
- Functional programs represent deep architecture
- **Symbolic program synthesis** to choose:
 - Which neural library functions for re-use
 - Architecture that puts them together

Houdini 

Houdini



Synthesis of Differentiable Functional Programs



Language

Types

guides the search space

$$\begin{aligned}\tau & ::= \textit{Atom} \mid \textit{ADT} \mid F \\ \textit{Atom} & ::= \texttt{bool} \mid \texttt{real} \\ \textit{TT} & ::= \textit{Atom} \mid \textbf{Tensor}\langle\textit{Atom}\rangle[m_1][m_2] \dots [m_k] \\ \textit{ADT} & ::= \textit{TT} \mid \alpha\langle\textit{TT}\rangle \\ F & ::= \textit{ADT} \mid F_1 \rightarrow F_2.\end{aligned}$$

Programs

defines the search space

$$e ::= \langle\!\langle \tau \rangle\!\rangle \mid \oplus_w \mid e_0 \circ e_1 \mid \mathbf{map}_\alpha e \mid \mathbf{fold}_\alpha e \mid \mathbf{conv}_\alpha e.$$

↓

↓

“neural functions”, e.g., two-layer ff network

list and graph versions

Search

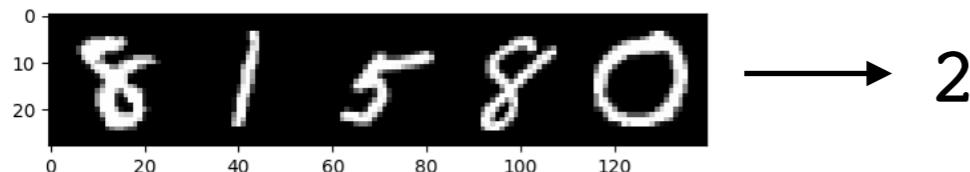
Enumeration of programs in language

Shortest to longest

Types reduce search space

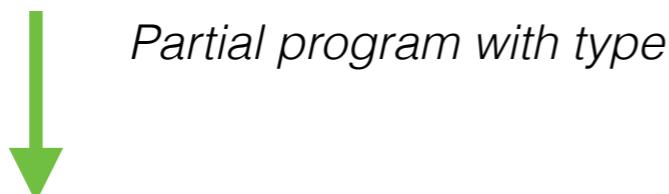
Limit number of trainable functions per candidate

Training set



Search

<< Tensor[28,148] -> real >>



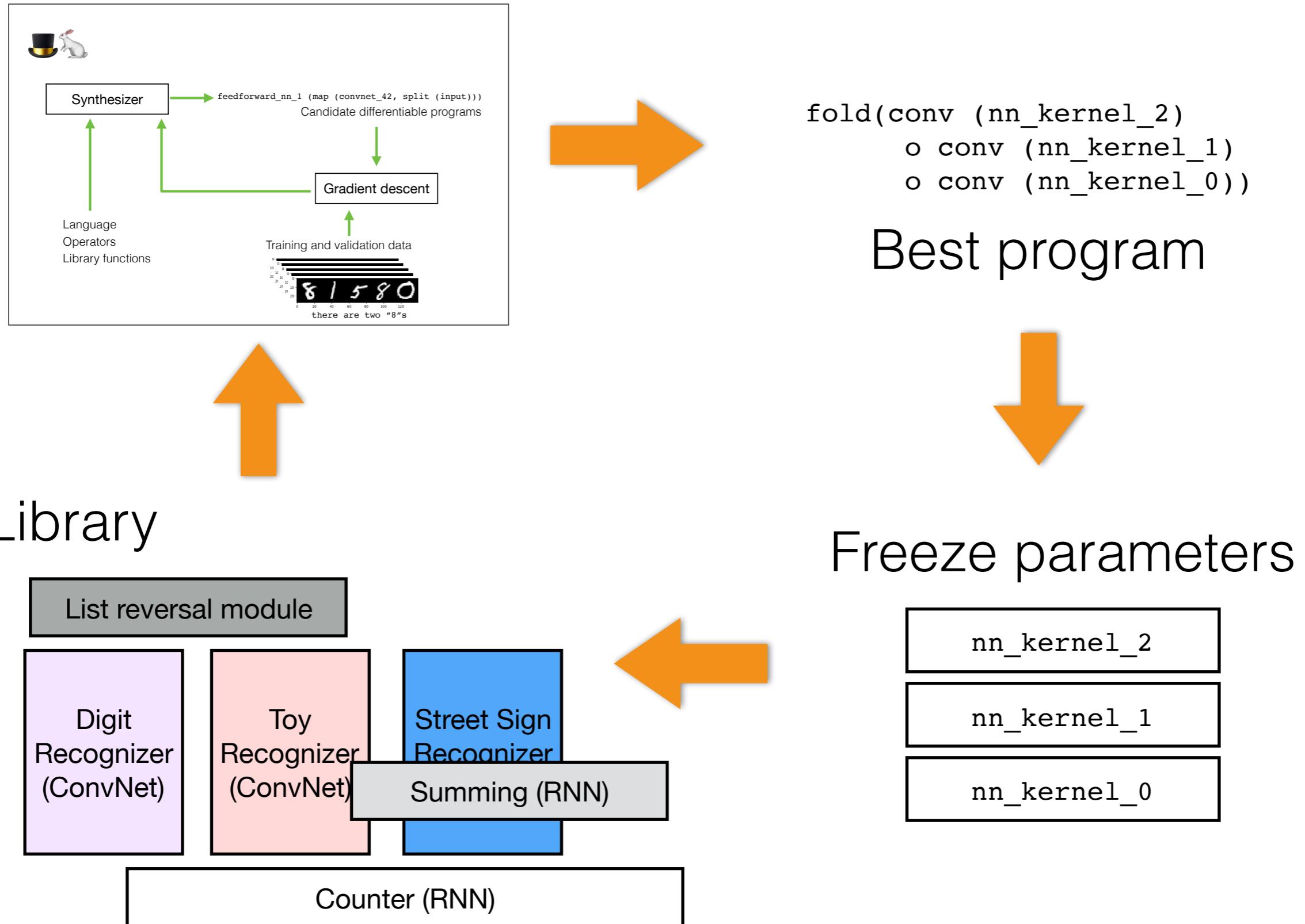
Expand using grammar

map

fold[list] << Tensor[28,148] -> list[T] >>



Synthesis for Lifelong Learning

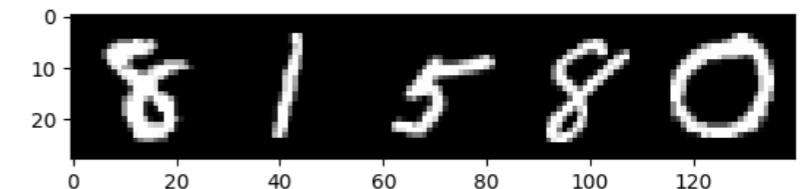


Experiments

Counting

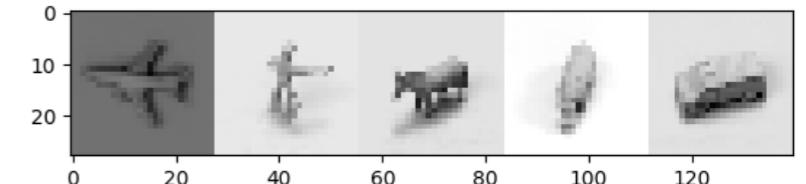
MNIST digits:

Recognize(5); Recognize(8); Count(8)



MNIST digits:

Recognize(5), Count(5); Count(8); Recognize(5)



MNIST/NORB:

Recognize(5), count(5), count(toy airplane), recognize(toy airplane)

Summing

Classify(1...10), Sum(sequence)

All-pairs shortest path

$(\text{conv}_{\text{graph}(\text{Tensor}[2])}^i nn_relax) \circ (\text{map}_{\text{graph}(\text{Tensor}[32][32][3])} perceive)$

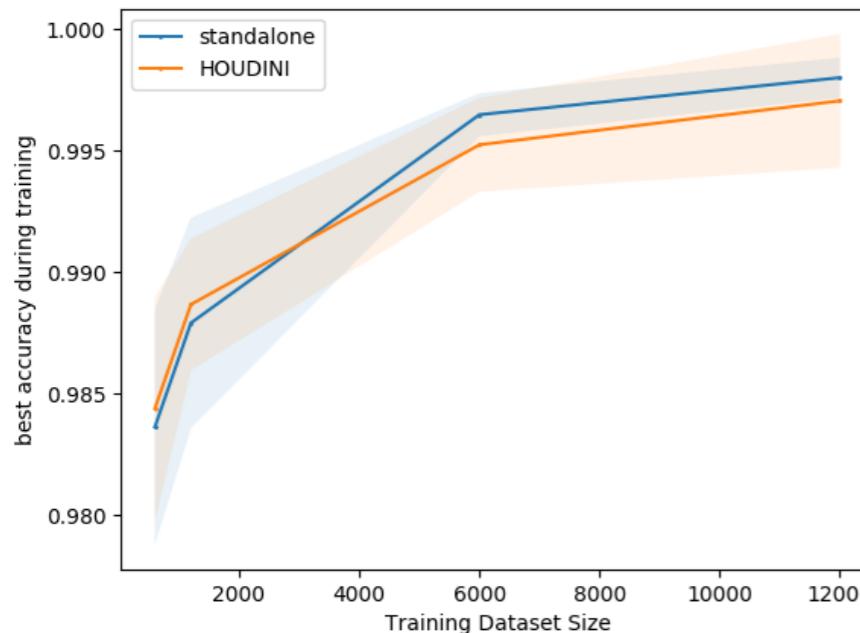


Classify(sign), Shortest_path(sign)

Classify(mnist), Shortest_path(mnist), Shortest_path(mnist)

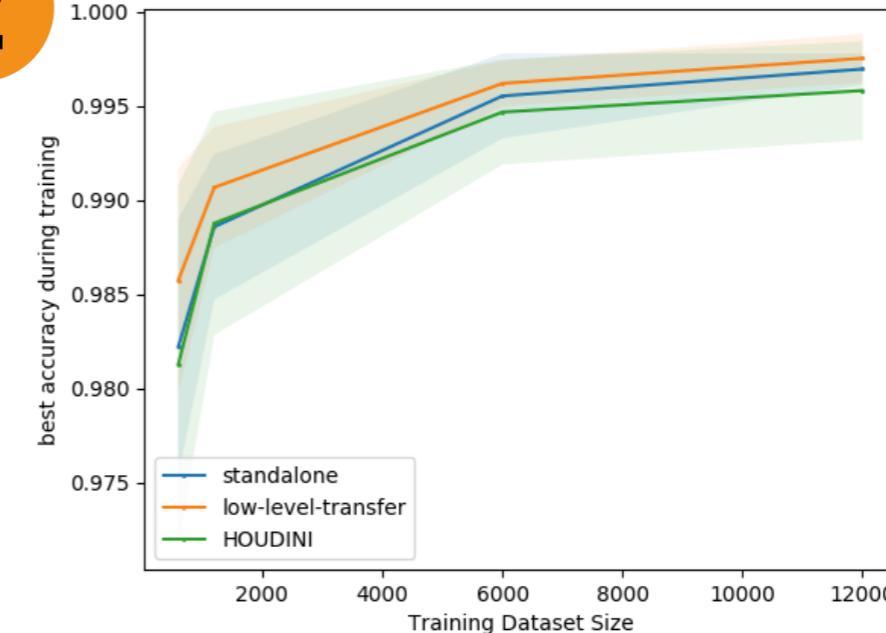
Results: Learning to Count

1



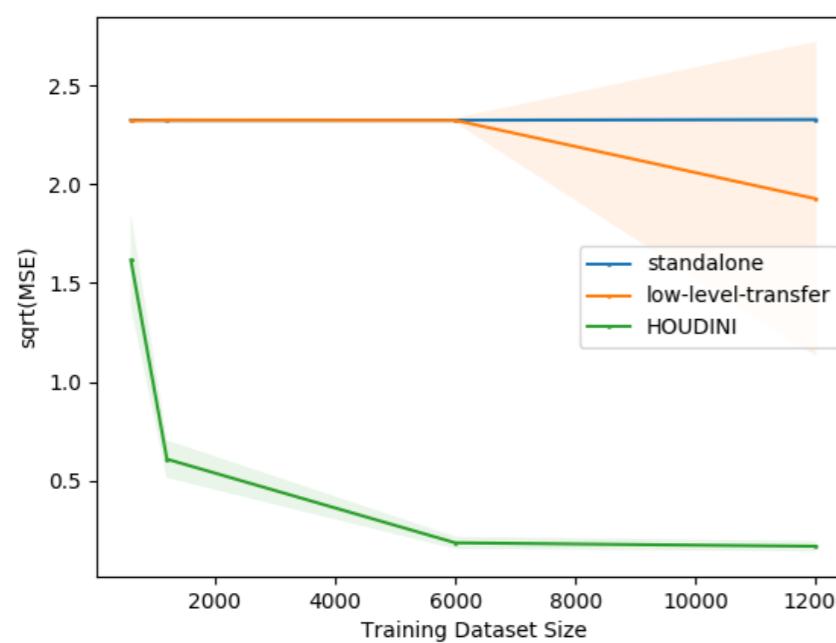
Recognise digit d_1
(binary classification)

2



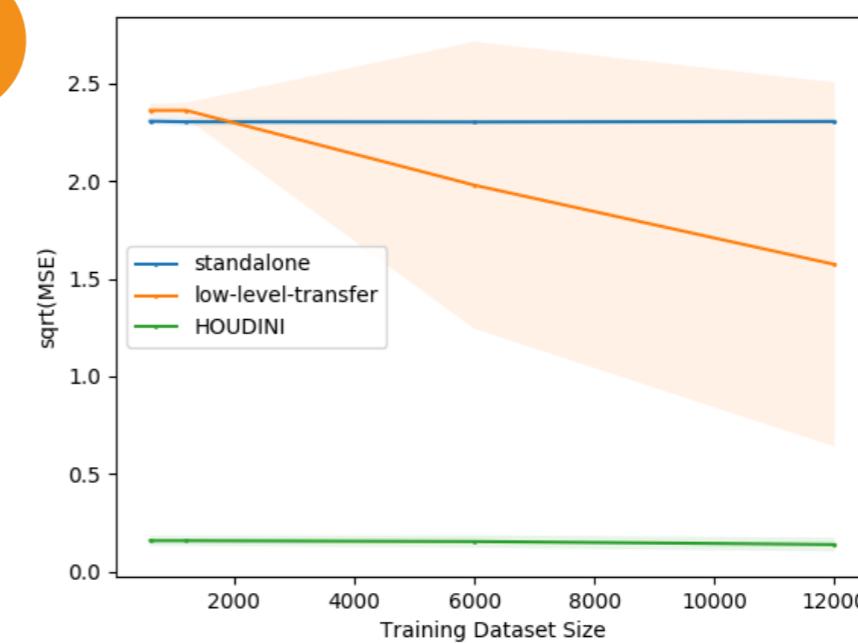
Recognise digit d_2
(binary classification)

3



Count digits d_1

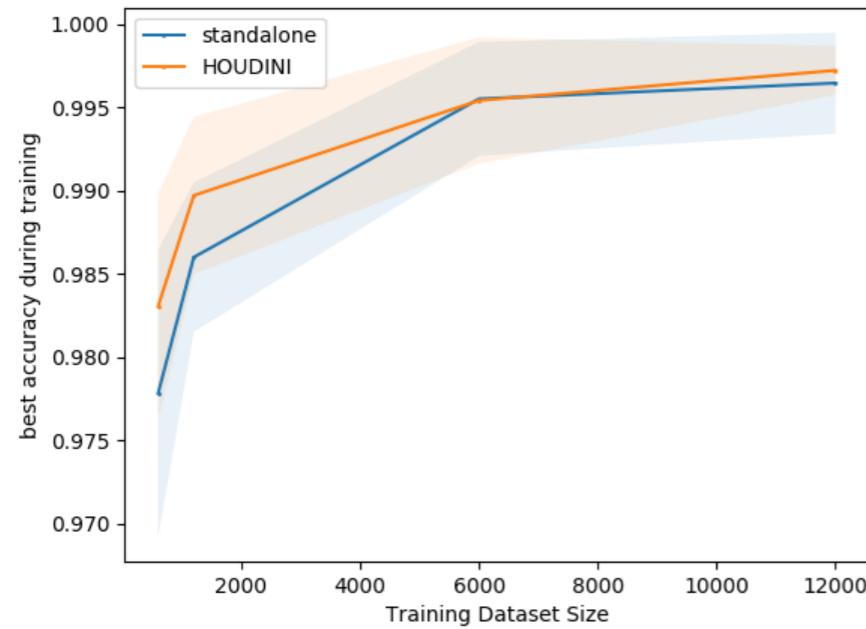
4



Count digits d_2

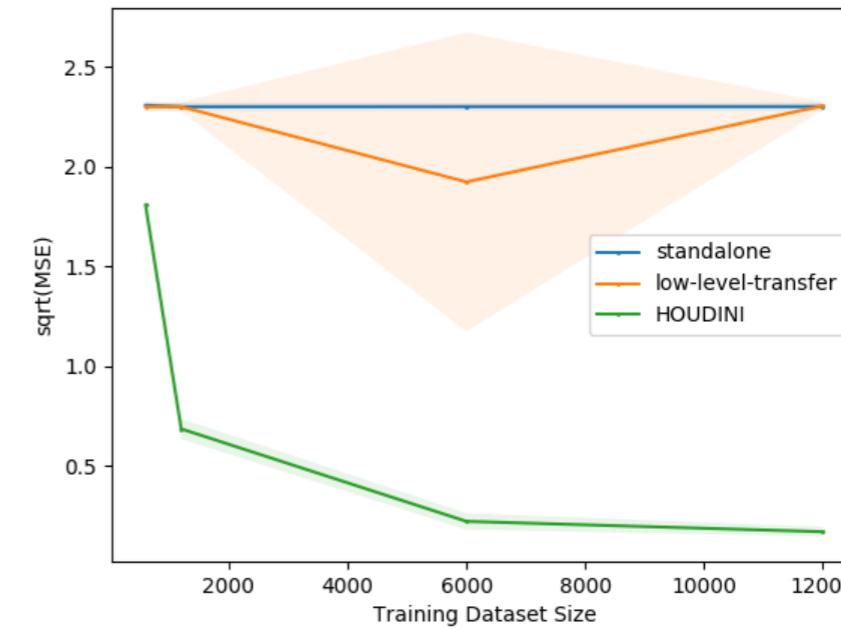
Results: Counting Toys

1



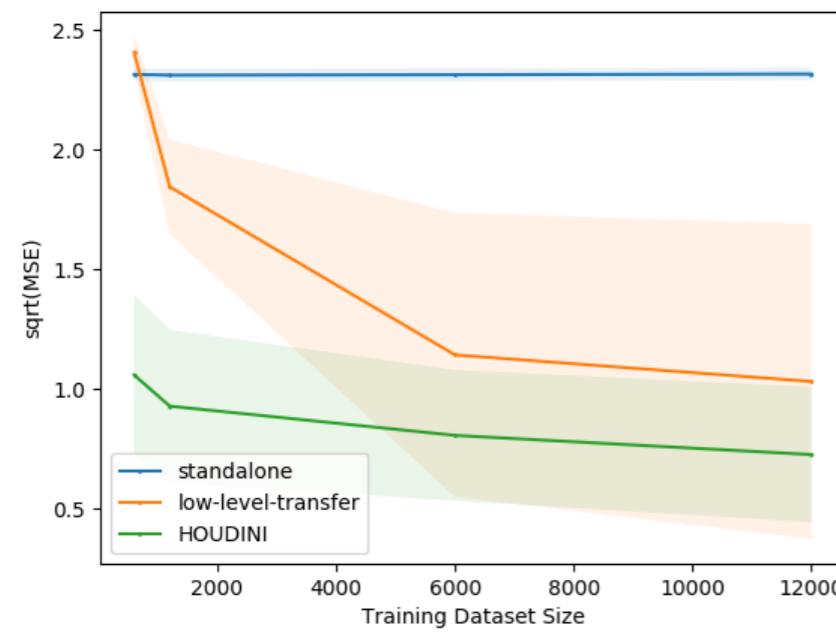
Recognise digit d_1

2



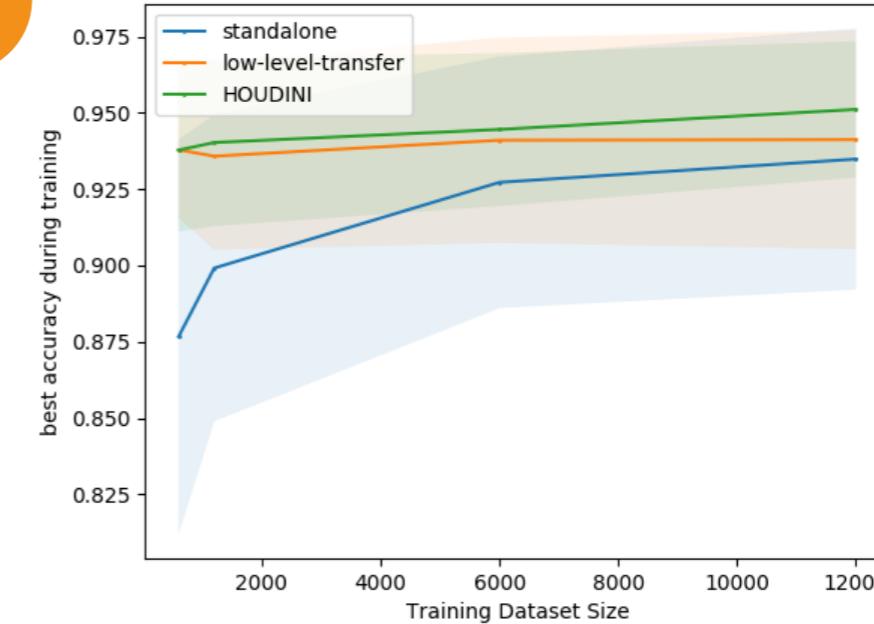
Count digit d_1

3



Count toys t_1

4



Recognise toy t_1

Results: Shortest path

Training: Grids of size 2x2, 3x3, 4x4

Testing: Grids of size 5x5

| | Image → Node cost (RMSE) | Shortest path length (RMSE) |
|-----------------------|-----------------------------|--------------------------------|
| Vanilla CNN → LSTM | 0.37 | 5.97 |
| Houdini | 0.38 | 1.53 |

Street sign images *Street sign images*

Results: Shortest path (transfer)

Training: Grids of size 2x2, 3x3, 4x4

Testing: Grids of size 5x5

| | Image → Node cost (RMSE) | Shortest path length (RMSE) | Shortest path length (RMSE) |
|-----------------------|-----------------------------|-----------------------------------|-----------------------------------|
| Vanilla CNN → LSTM | 1.21 | 5.33 | 6.16 |
| Houdini | 1.29 | 1.62 | 4.98 |

MNIST images *MNIST images* *Street sign images*

Impact of type system

| Program depth | 4 | 5 | 6 |
|----------------|-------|--------|---------|
| No type system | 15633 | 247589 | 3449845 |
| Type system | 25 | 155 | 444 |

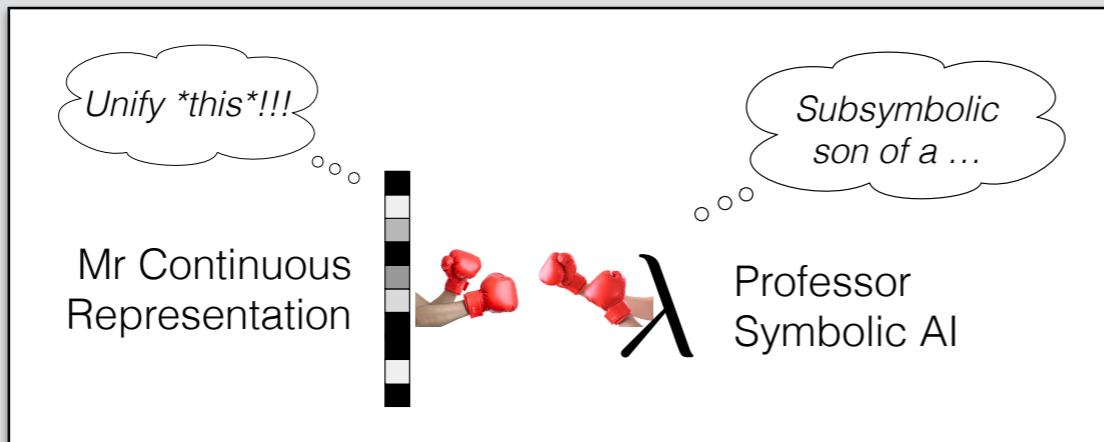
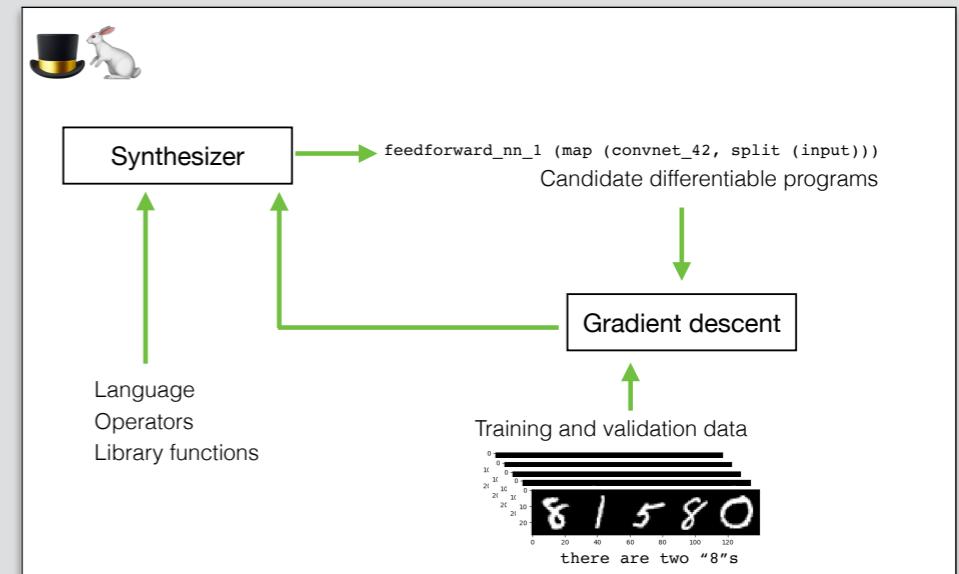
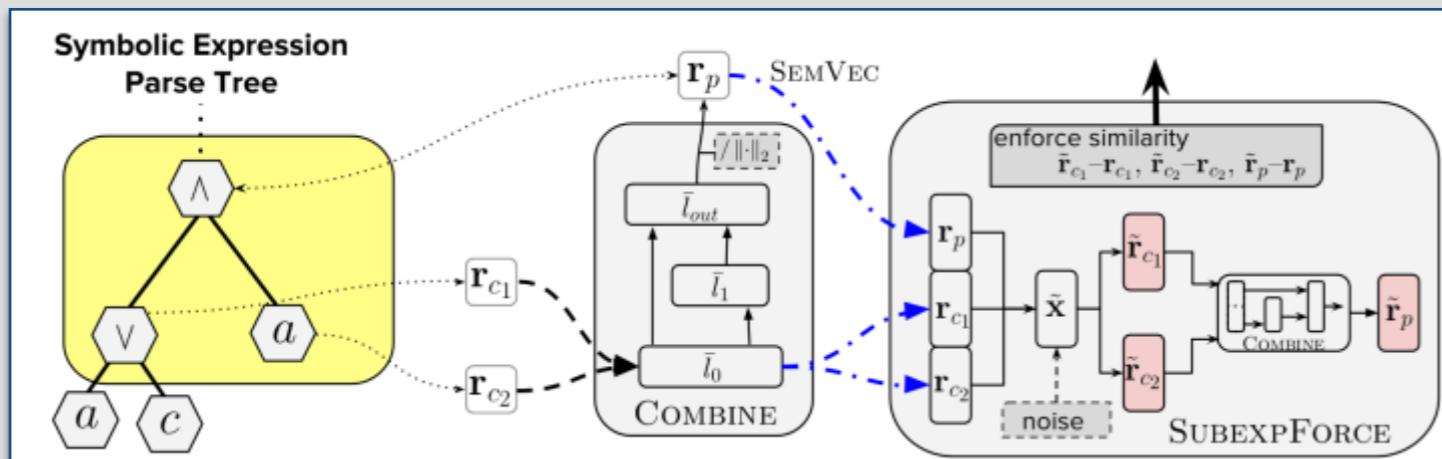
Types of transfer

- Low-level transfer
 - Reuse perceptual network across high-level tasks
- High-level transfer
- “Infilling”
 - Learn perceptual concepts given only supervision at high level
- Selective transfer
 - Synthesis algorithm decides whether and when to re-use

Adventures in Neurosymbolic Learning

Charles Sutton

bit.ly/adventures-neurosymbolic



Thanks!

Miltos Allamanis
Pushmeet Kohli
Pankajan Chanthirasegaran
Kai Xu
Lazar Valkov
Dipak Chaudhari
Akash Srivastava
Swarat Chaudhuri