

Learning Program Representations: Symbols to Vectors to Semantics

Charles Sutton
University of Edinburgh
& The Alan Turing Institute

10 December 2016

<http://edin.ac/2ggR9uK>

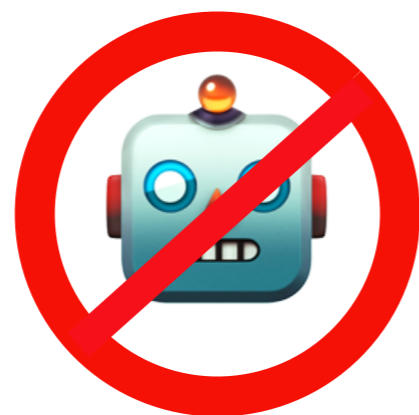


THE UNIVERSITY of EDINBURGH
informatics

**The
Alan Turing
Institute**

Microsoft
Research

EPSRC
Engineering and Physical Sciences
Research Council



Source code is a means of human communication

```
try{  
    Node $name=$methodInvoc();  
    $BODY$  
}finally{  
    $(Transaction).finish();  
}
```

```
public static final  
String $name = $StringLit;
```



```

ConfigurationBuilder.<init>
ConfigurationBuilder.setOAuthConsumerKey
ConfigurationBuilder.setOAuthConsumerSecret
ConfigurationBuilder.build
TwitterFactory.<init>
TwitterFactory.getInstance

```

```

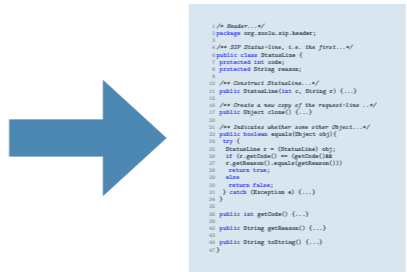
try{
    Node $name=$methodInvoc();
    $BODY$
}finally{
    $(Transaction).finish();
}

```

```

1 /* Header */
2 package org.zoolu.sip.header;
3
4 /** SIP Status-Line, i.e. the first
5  * line of a response message */
6 public class StatusLine {
7     protected int code;
8     protected String reason;
9
10    /** Construct StatusLine */
11    public StatusLine(int c, String r) {
12        code = c;
13        reason = r;
14    }
15
16    /** Create a new copy of the request-line */
17    public Object clone() {
18        return new StatusLine(getCode(), getReason());
19    }
20
21    /** Indicates whether some other Object
22     * is "equal to" this StatusLine */
23    public boolean equals(Object obj){
24        try {
25            StatusLine r = (StatusLine) obj;
26            if (r.getCode() == getCode() &&
27                r.getReason().equals(getReason()))
28                return true;
29            else
30                return false;
31        } catch (Exception e) {
32            return false;
33        }
34    }
35
36    public int getCode() {
37        return code;
38    }
39
40    public String getReason() {
41        return reason;
42    }
43
44    public String toString() {
45        return "SIP/2.0 " + code + " " + reason + "\r\n";
46    }
47 }

```



Summarisation

Topic models [ICSE 2016]

Learning how libraries are used

Nonparametric Bayes grammars [FSE 2014]

Probabilistic pattern mining [KDD 2016; FSE 2016]



Learning coding conventions

[FSE 2014, 2015]

-
- Further ahead...
 - Defining requirements
 - Architecting
 - Navigation
 - Maintenance
 - Optimising performance
 - Testing, verification
 - Refactoring
 - Porting
 - Debugging
-

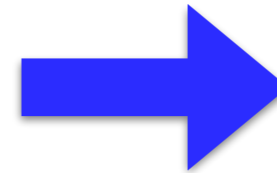
Neural networks that capture program semantics

Why not solved?

1. Domain connections
2. Typical programs

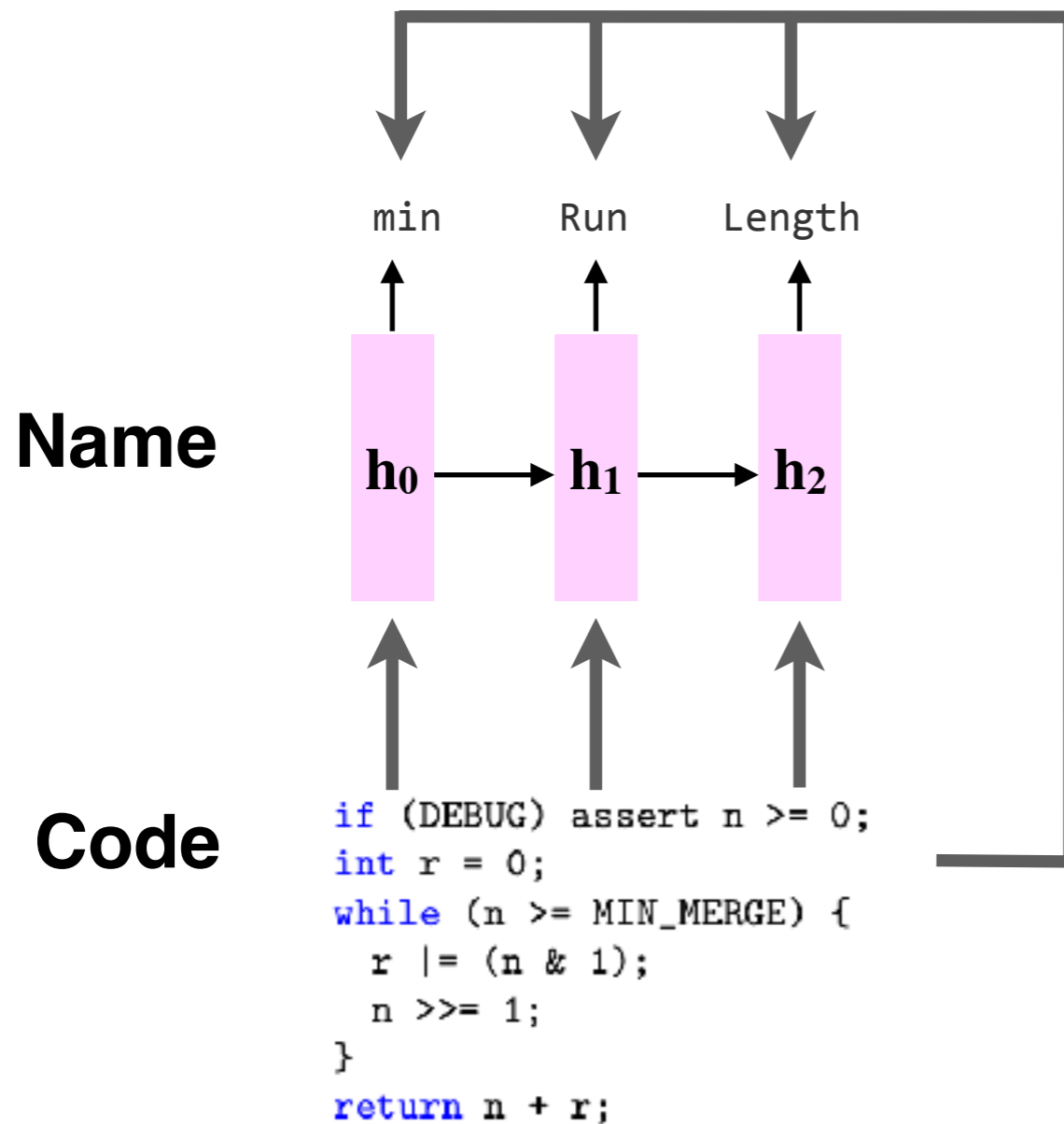
Learning to Name

```
if (DEBUG) assert n >= 0;
int r = 0;
while (n >= MIN_MERGE) {
    r |= (n & 1);
    n >>= 1;
}
return n + r;
```

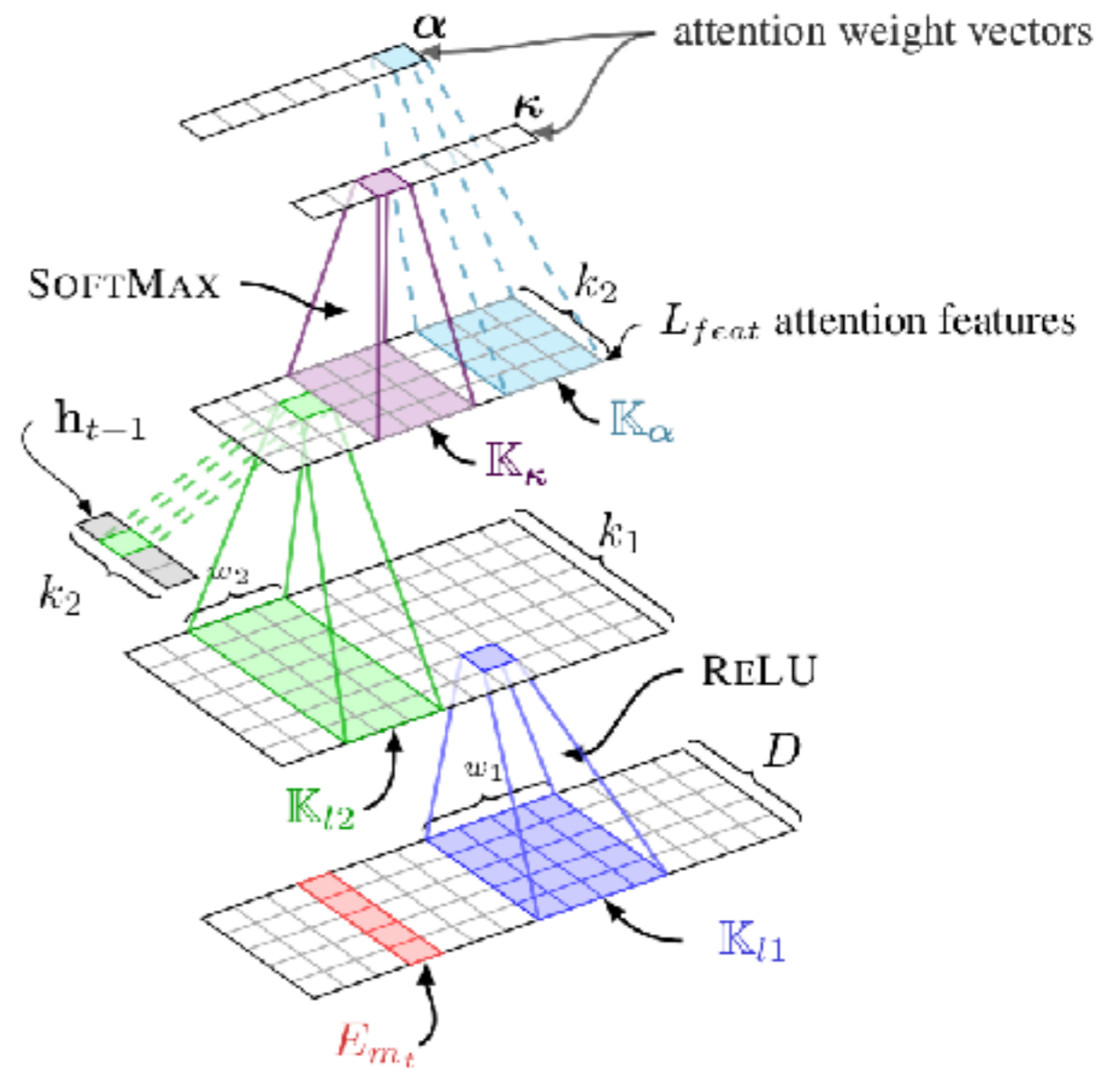


minRunLength

Predicting Names of Methods



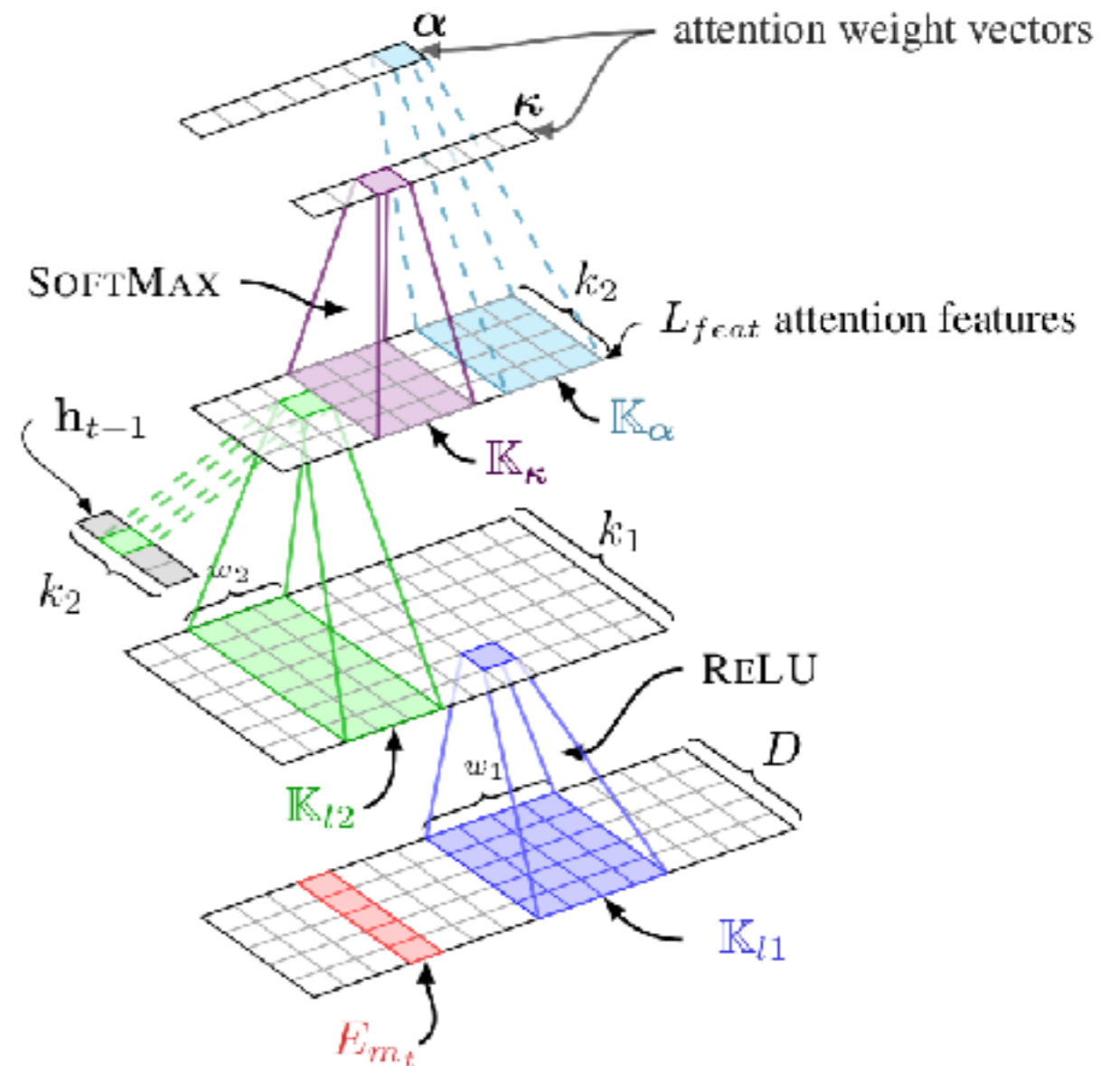
RNN for generating summary



convolutional attention mechanism

Three Attention Mechanisms

- α : Distribution over input locations
 - Weights for averaging input embeddings
- κ : Distribution over input locations
 - Weights for copying tokens from input to output (even OOV)
 - Related to pointer networks
[Vinyals et al, 2015]
- λ : Scalar $[0, 1]$
 - weight to decide two mechanisms



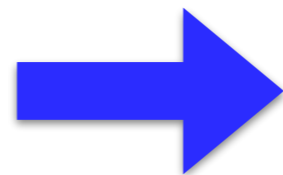
	F1 (%)		Exact Match (%)	
	Rank 1	Rank 5	Rank 1	Rank 5
tf-idf	40.0	52.1	24.3	29.3
Standard Attention	33.6	45.2	17.4	24.9
conv_attention	43.6	57.7	20.6	29.8
copy_attention	44.7	59.6	23.5	33.7

Standard attention: [Bahdanau, Cho, and Bengio, 2015]

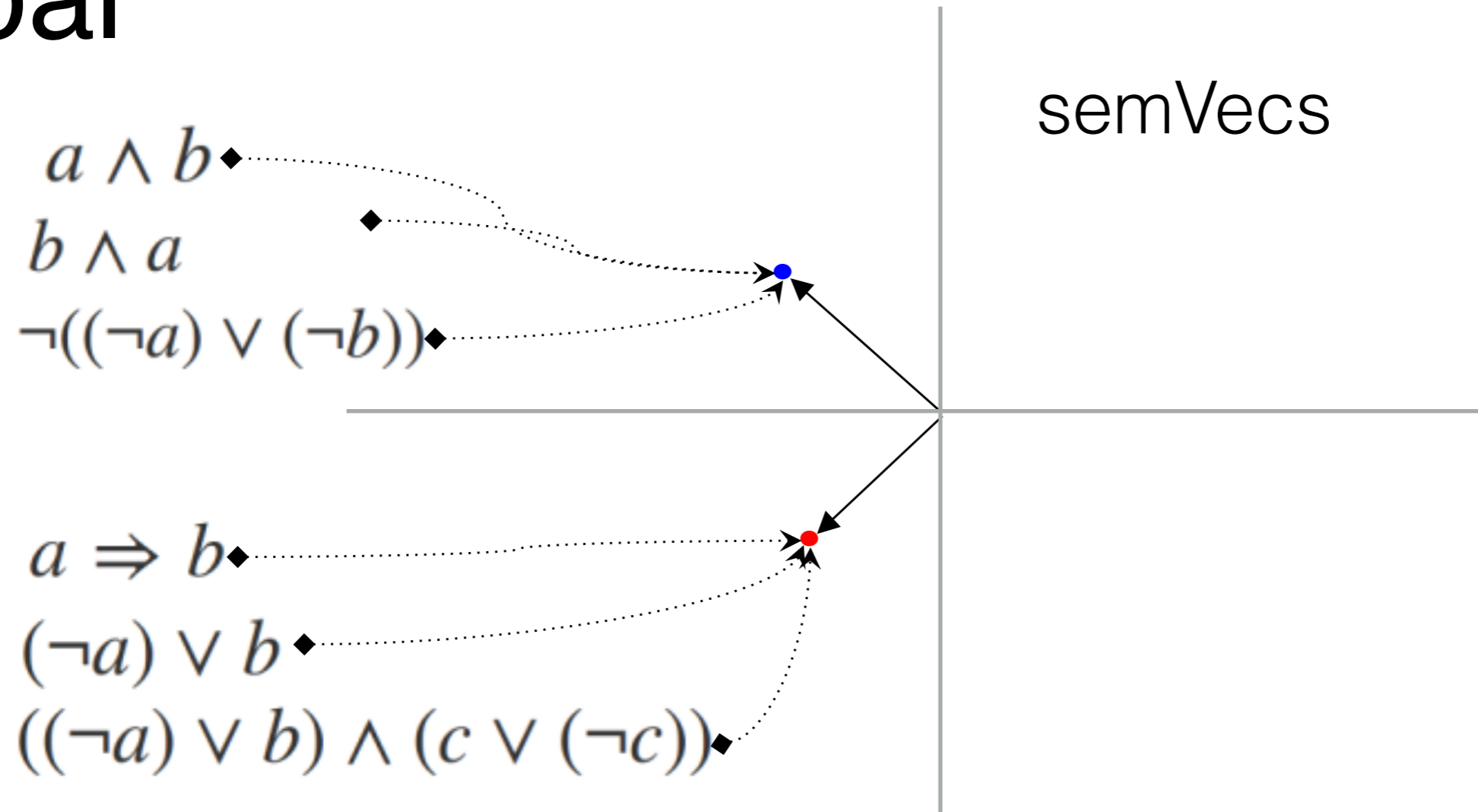
Continuous Semantics for Symbolic Expressions

[Allamanis, Chanthirasegaran Kohli, and Sutton, 2017]

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



Goal



How much symbolic semantics (semantic equivalence)

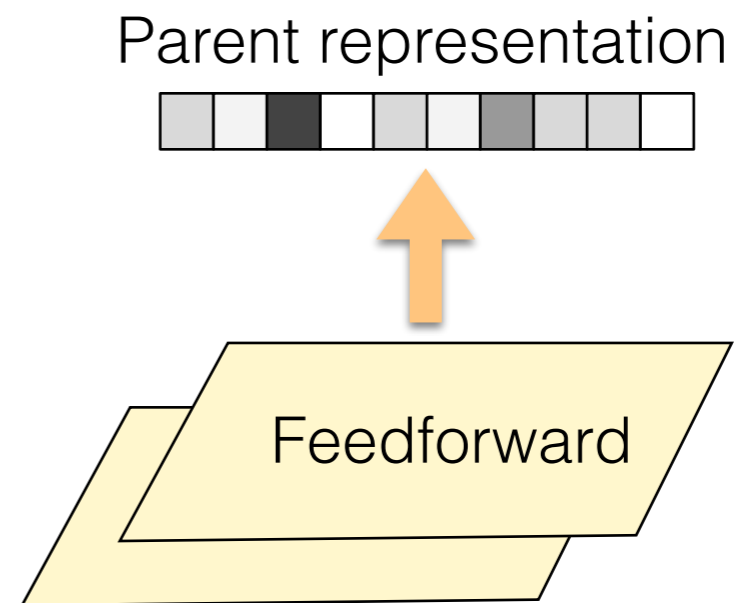
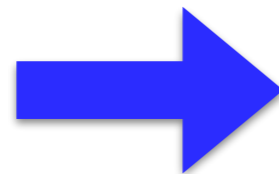
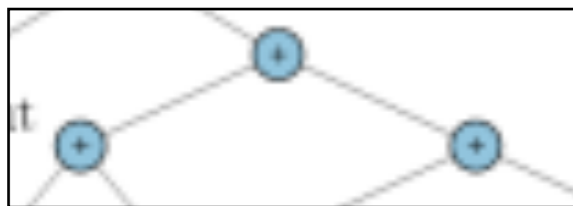
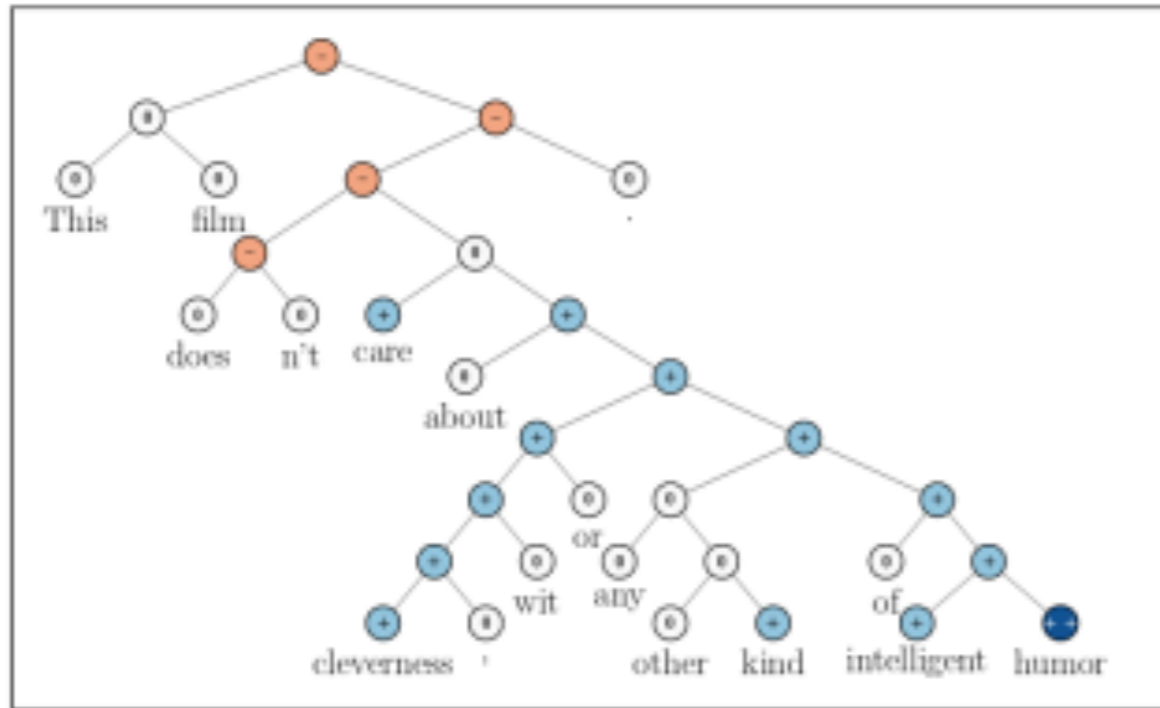
Assume we compress into **continuous** vector?

Symbolic reasoning:

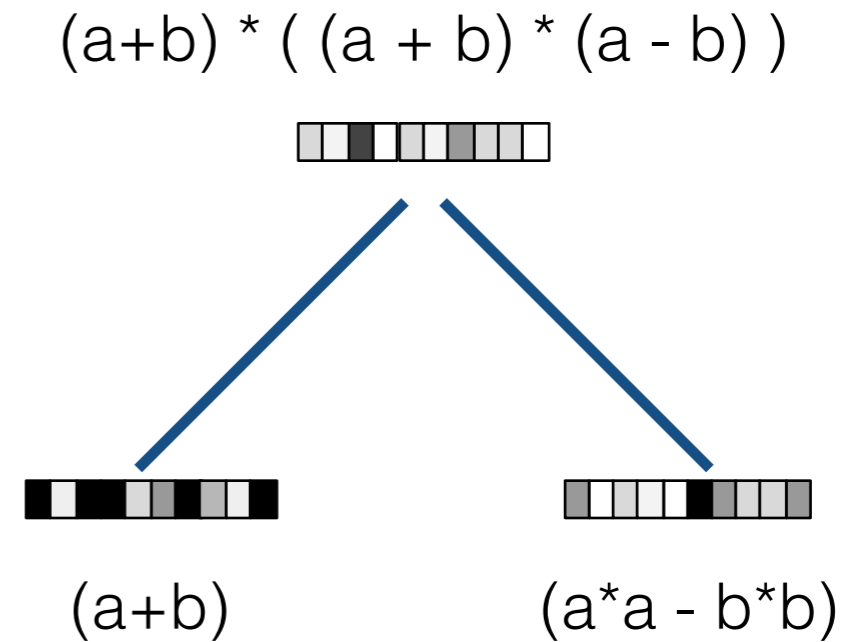
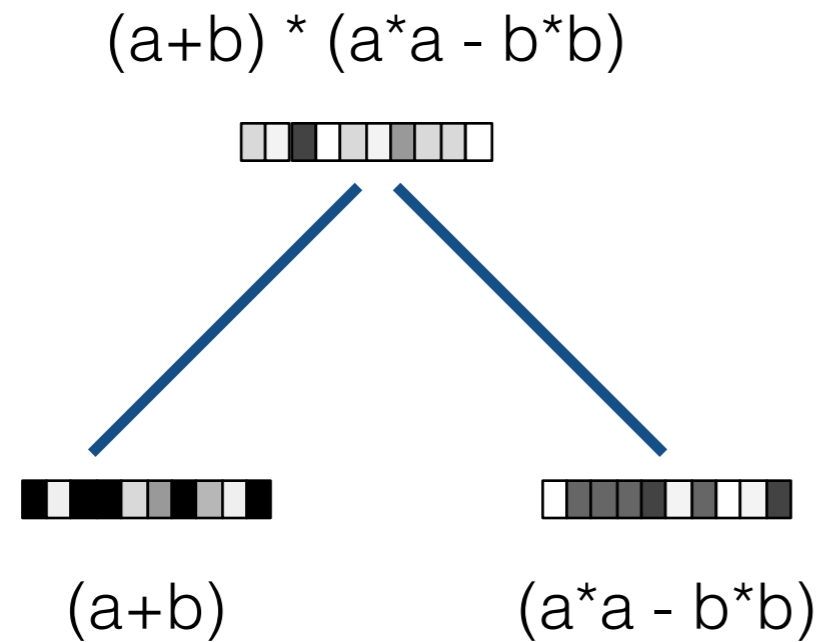
~~search~~

pattern recognition

Recursive NN (TreeNN)



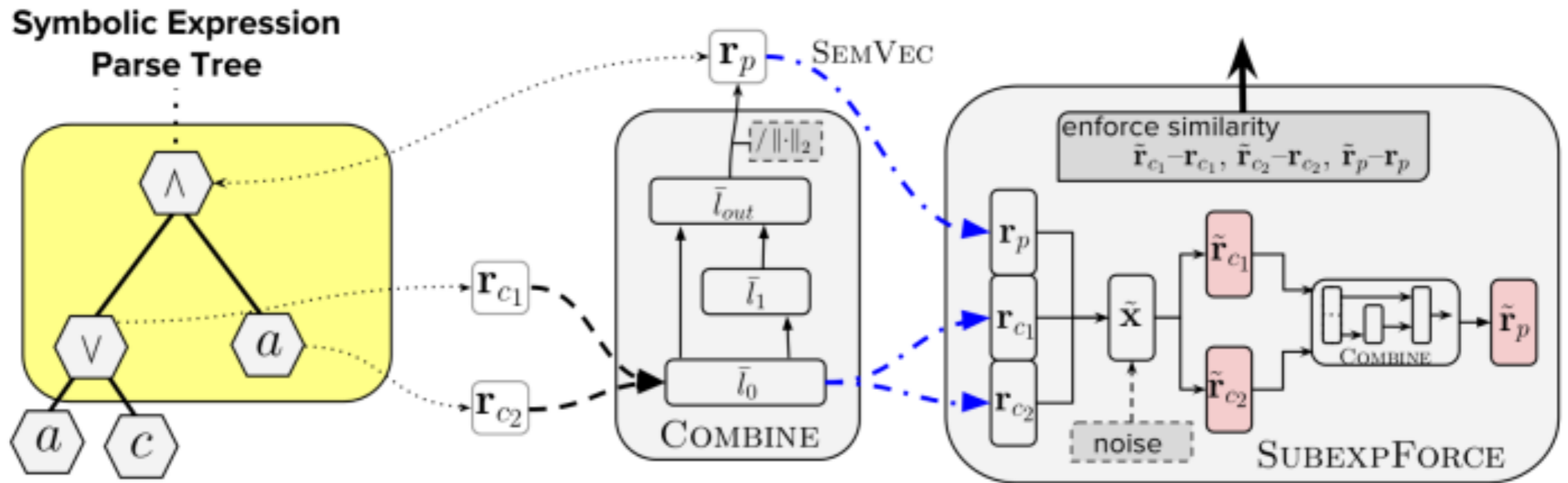
Problem: Separating out syntax



semantically equivalent, different vectors!

Result: nearest neighbours mostly reflect syntax

EqNet



Motivation via Unification

Semantic information is bidirectional

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

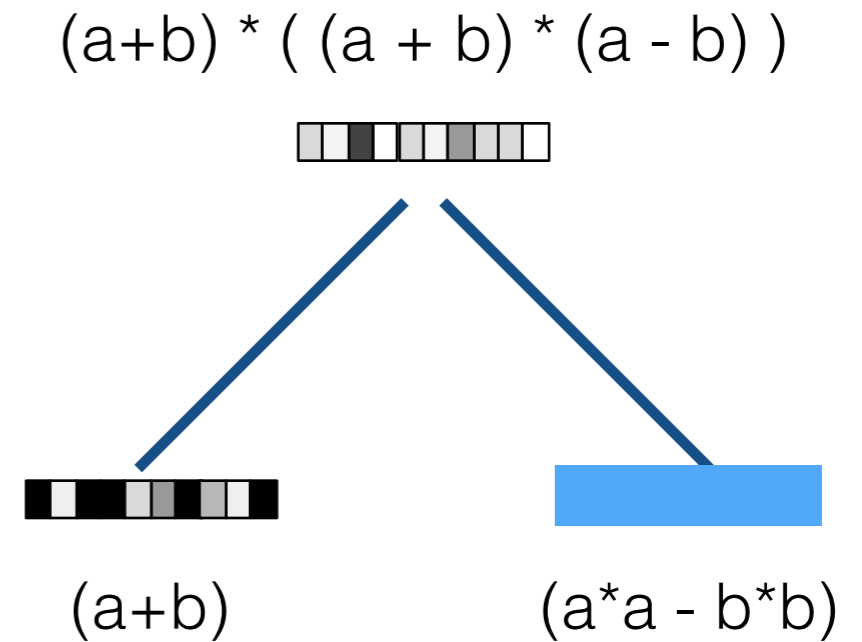
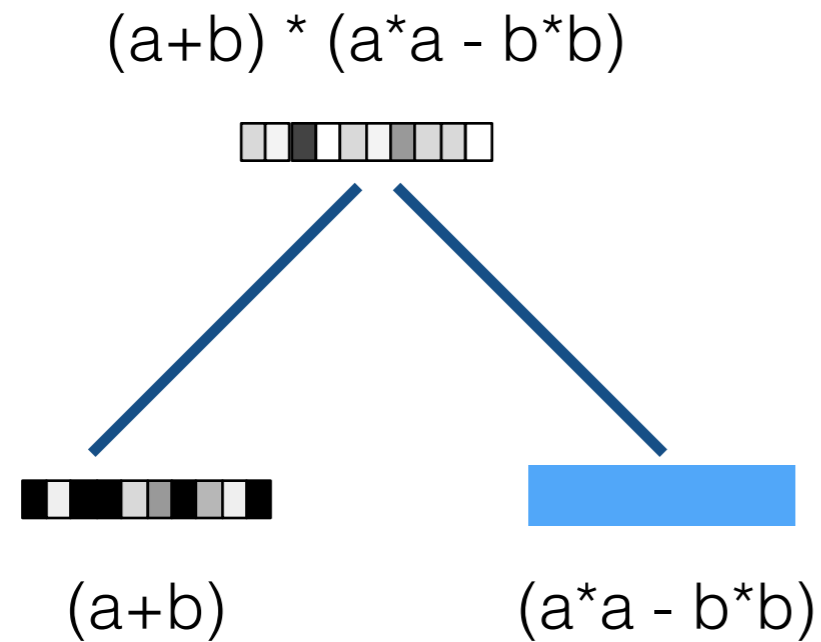
But **parents** provide info re **children**

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

Unification propagates this info automatically

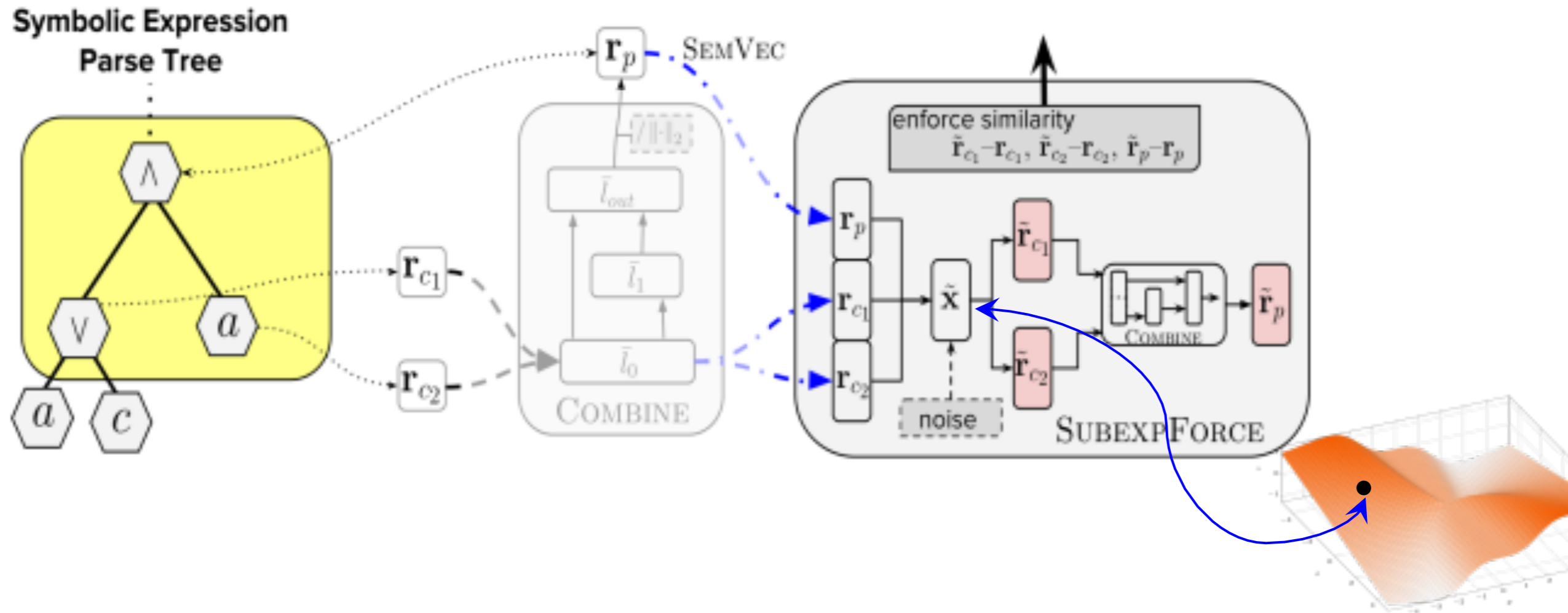
How to map to continuous space?

Subexpression Forcing



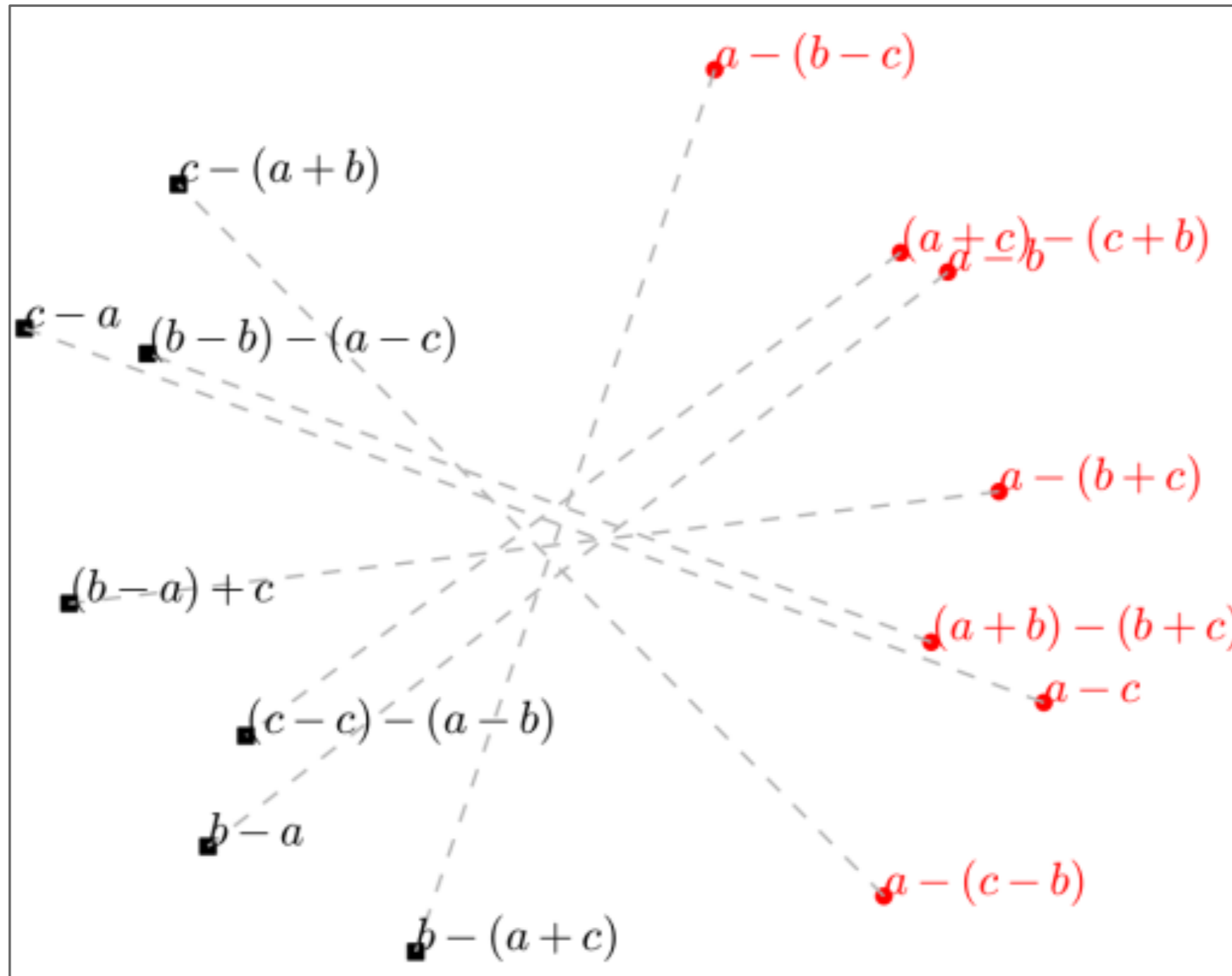
ensure this prediction problem is “easy”
semantic classes will be clustered together

Subexpression Forcing



Denoising autoencoder plus bottleneck
on (parent, child1, child2) representations
(Additional regulariser)

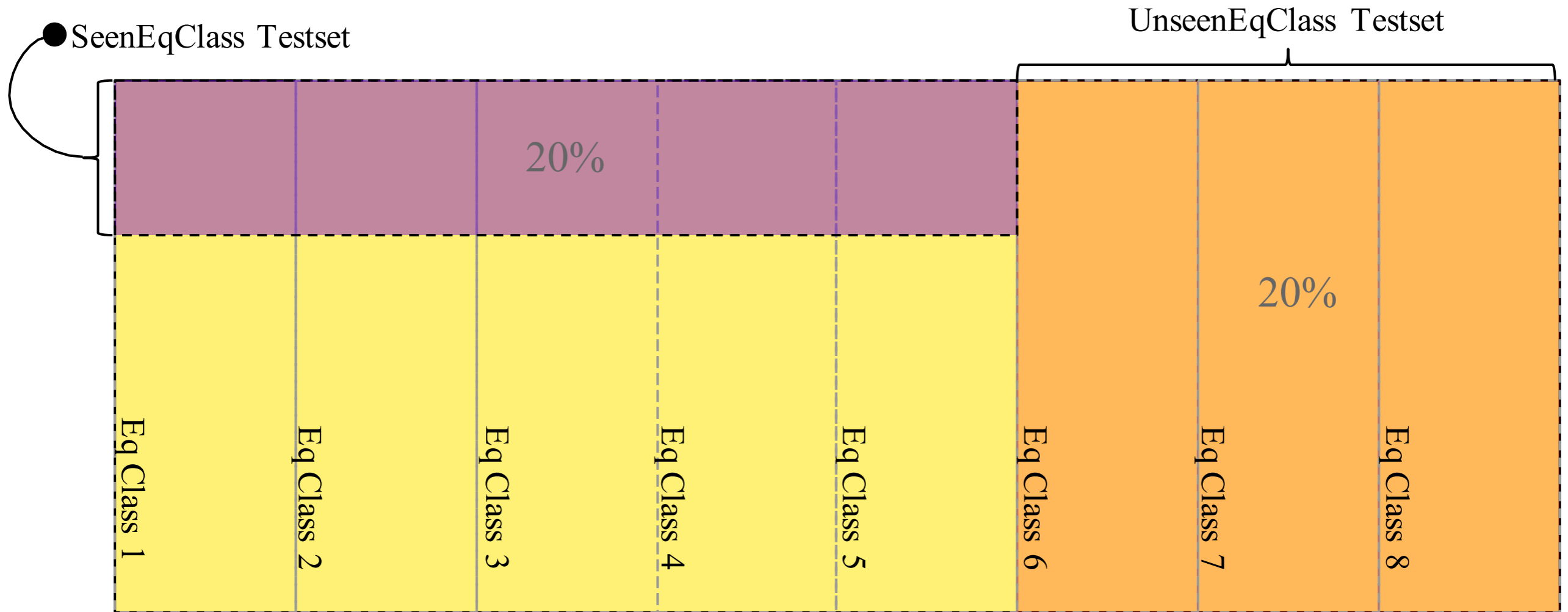
Visualizing polynomials



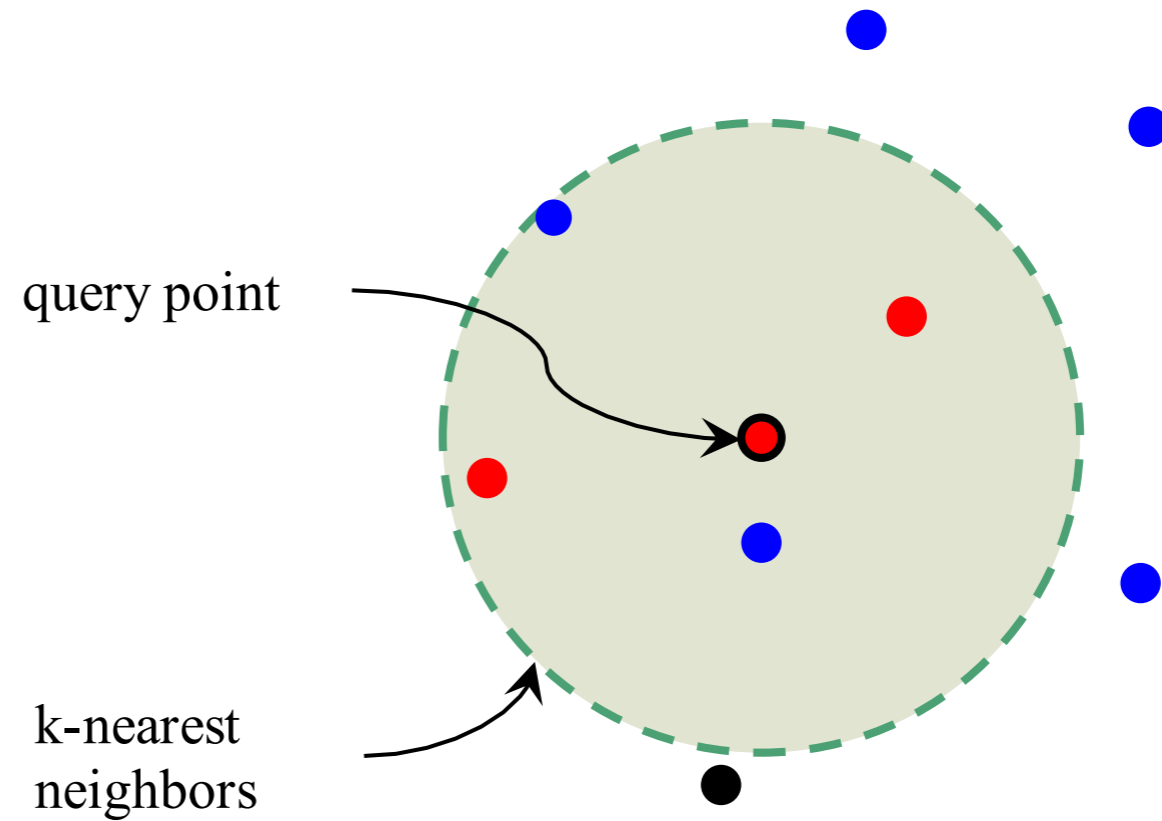
Evaluation

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

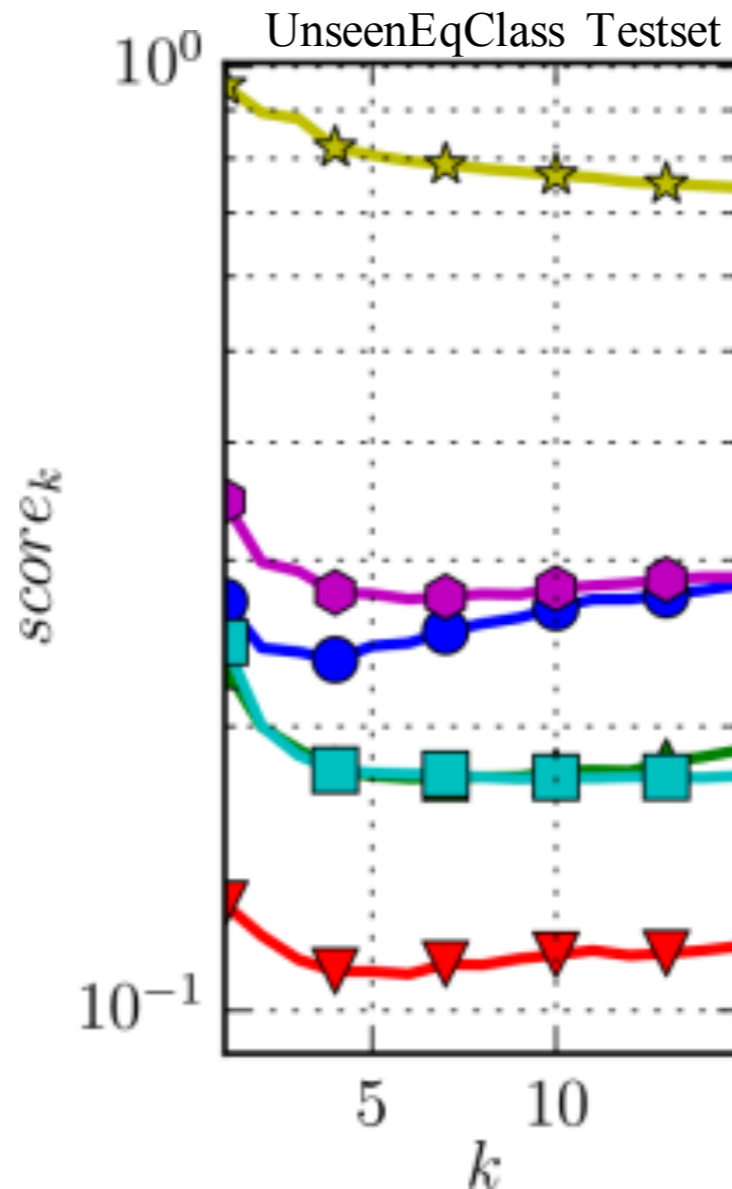
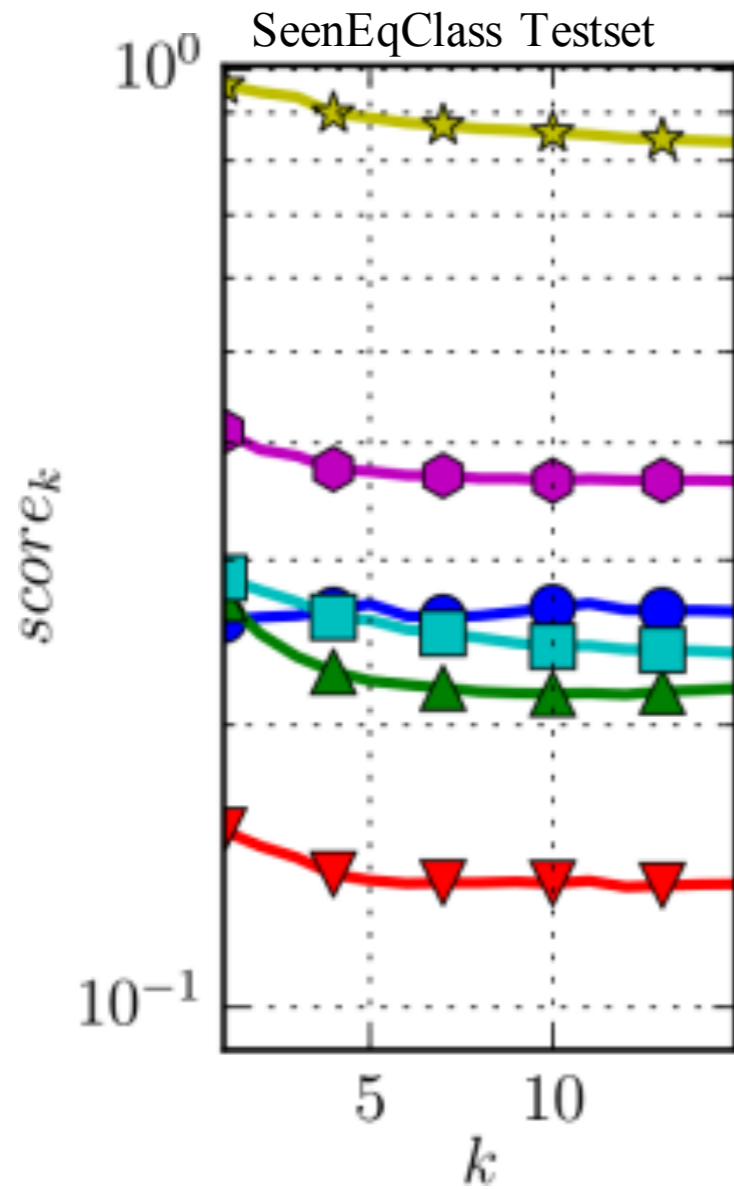


Evaluation Metric

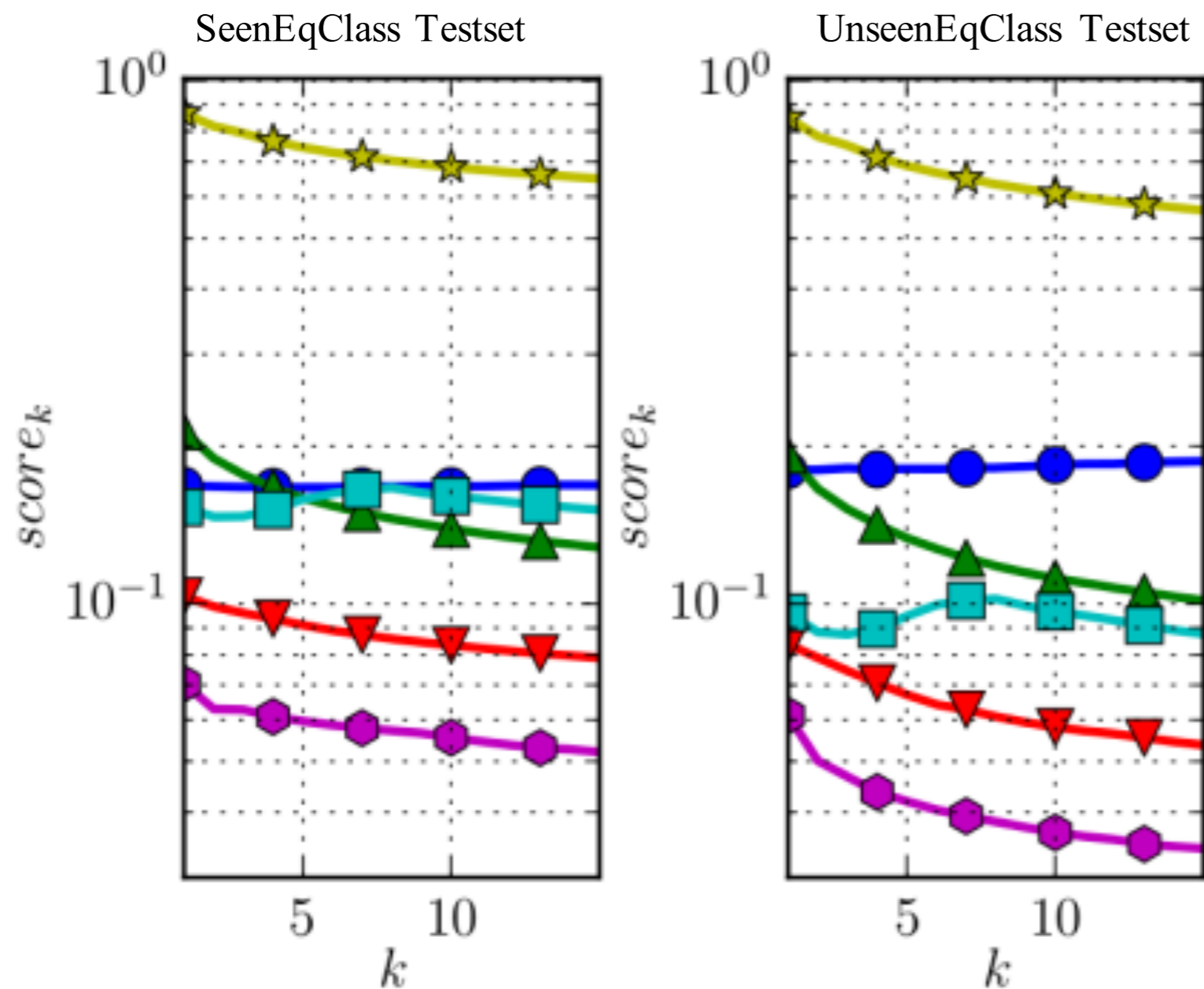


$$score_k(q) = \frac{|\mathbf{N}_k(q) \cap c|}{\min(k, |c|)}$$

Results

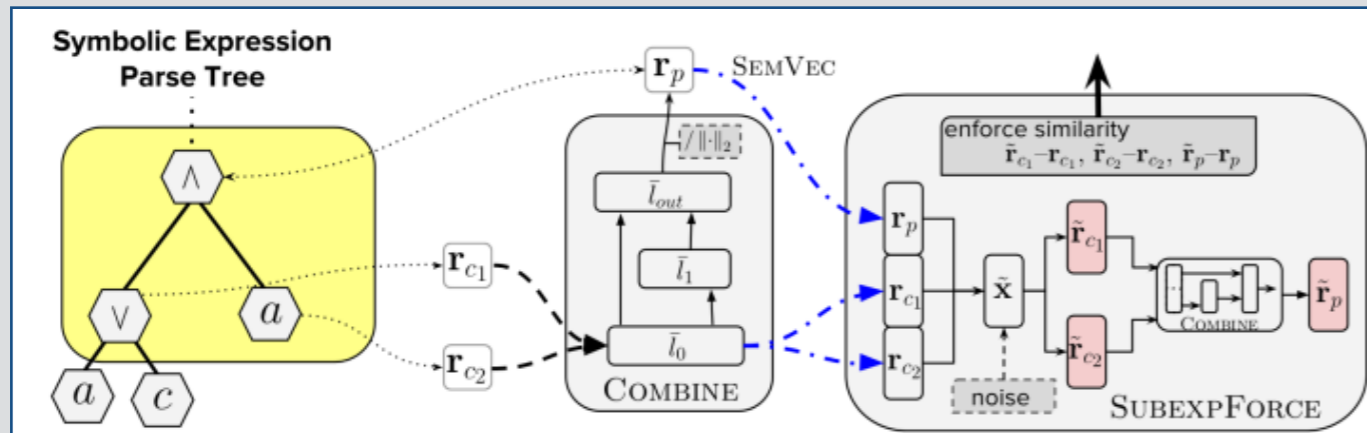


Evaluating compositionality

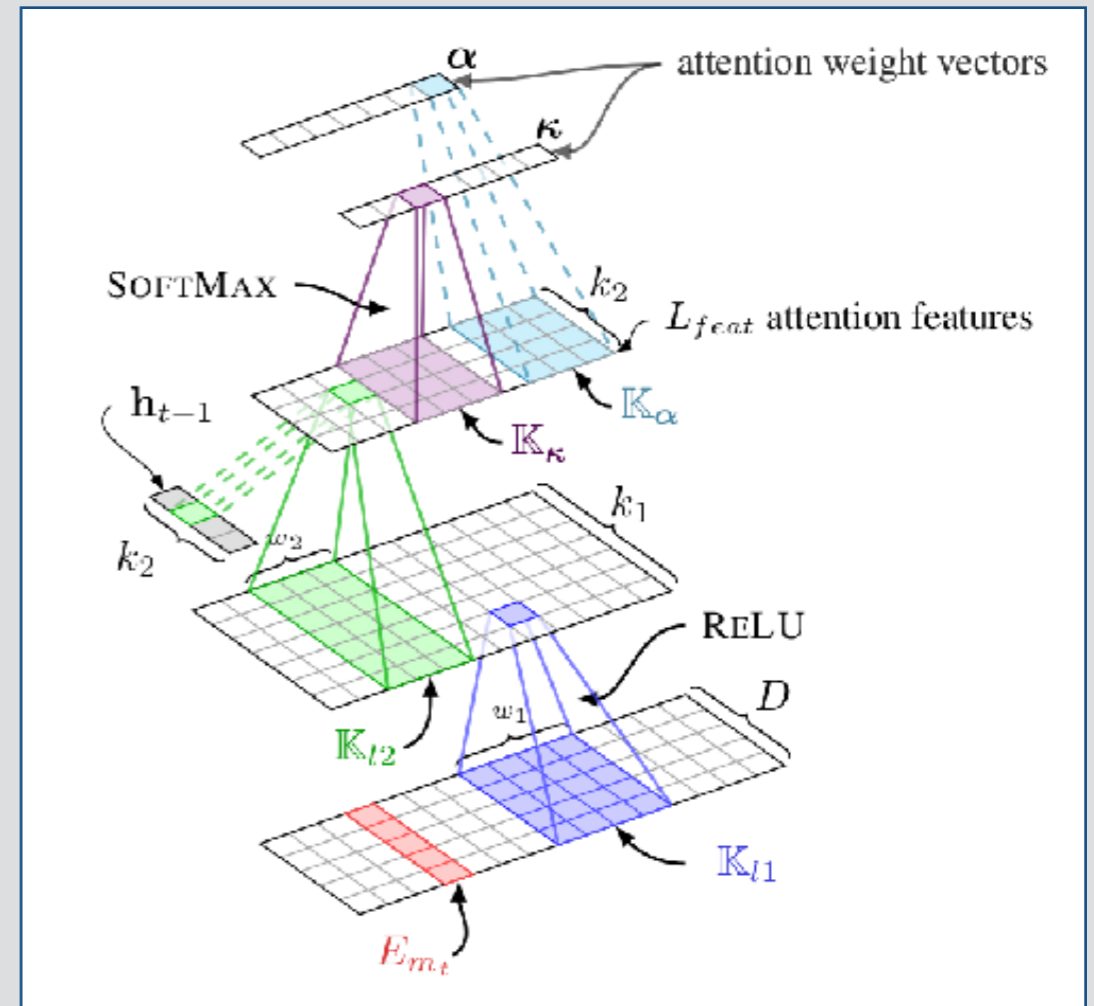


Learning Program Representations: Symbols to Vectors to Semantics

Charles Sutton, University of Edinburgh



Equivalence networks
for continuous semantics



Naming methods
convolutional attention

Thanks!

- Miltiadis Allamanis
- Hao Peng
- Pushmeet Kohli
- Pankajan Chantirasagar