

Experiments in Learning to Solve Formal Analogical Equations

Rafik Rhouma and Philippe Langlais



Dept. I.R.O.
Université de Montréal, Québec, Canada
`felipe@iro.umontreal.ca`

ICCBR 2018

Analogy

$[x : y :: z : t]$ stands for a **proportional analogy**

- ▶ x is to y as z is to t



Analogy

$[x : y :: z : t]$ stands for a **proportional analogy**

- ▶ x is to y as z is to t



Formal Analogies (or Analogies on Forms)

- ▶ $[x \setminus y \circ z \supset t]$ denotes a formal analogy

1 $[funny \setminus funniest \circ lucky \supset luckiest]$

2 $[sugg\grave{e}re \setminus sugg\grave{e}rer \circ ing\grave{e}re \supset ing\grave{e}rer]$

3 $[This\ guy\ drinks\ too\ much \setminus This\ boat\ sinks \circ$
 $These\ guys\ drank\ too\ much \supset These\ boats\ sank]$

- ▶ Several operational definitions of formal analogy $[?, ?]$

Definition of [?]

$[x \setminus y \circ z \supset t]$ **iff** we can find d -**factorizations** f_X, f_Y, f_Z and f_t such that, $\forall i \in [1, d]$:

$$(f_Y^{(i)}, f_Z^{(i)}) \in \left\{ (f_X^{(i)}, f_t^{(i)}), (f_t^{(i)}, f_X^{(i)}) \right\}$$

- ▶ $[this\ guy\ drinks\ too\ much \setminus this\ boat\ sinks \circ$
 $these\ guys\ drank\ too\ much \supset these\ boats\ sank]$ because:

x	≡	this	guy	€	dr	inks	too much
y	≡	this	boat	€	s	inks	€
z	≡	these	guy	s	dr	ank	too much
t	≡	these	boat	s	s	ank	€

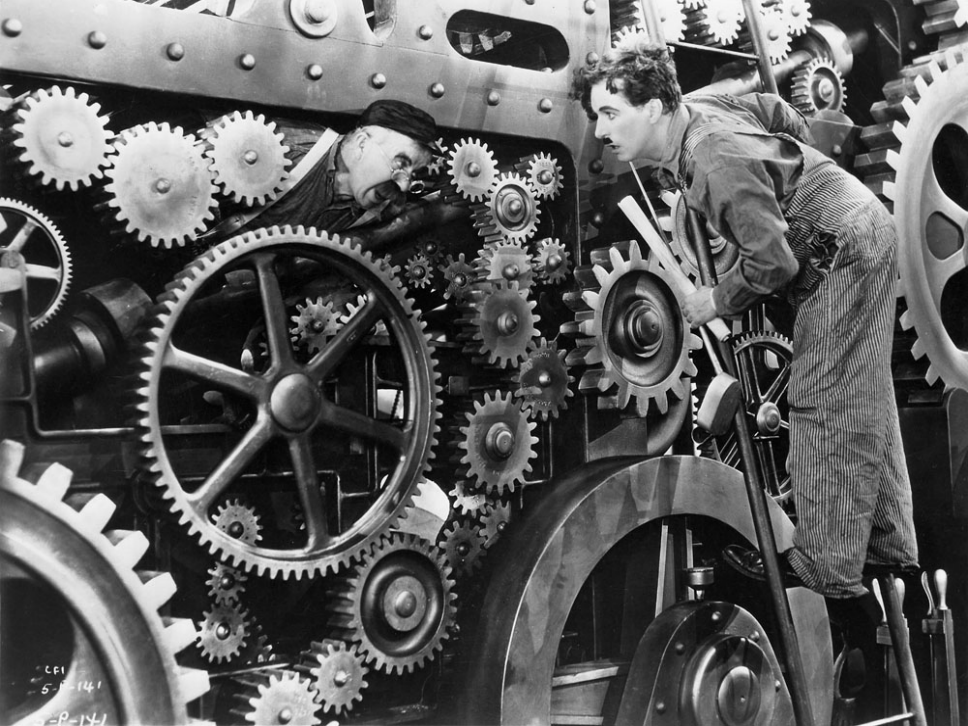
Definition of [?]

$[x \setminus y \circ z \supset t]$ **iff** we can find d -**factorizations** f_X, f_Y, f_Z and f_t such that, $\forall i \in [1, d]$:

$$(f_Y^{(i)}, f_Z^{(i)}) \in \left\{ (f_X^{(i)}, f_t^{(i)}), (f_t^{(i)}, f_X^{(i)}) \right\}$$

- ▶ $[$ *this guy drinks too much* \setminus *this boat sinks* \circ *these guys drank too much* \supset *these boats sank* $]$ because:

x	≡	this	guy	€	dr	inks	too much
y	≡	this	boat	€	s	inks	€
z	≡	these	guy	s	dr	ank	too much
t	≡	these	boat	s	s	ank	€



CFI
5-1-141

5-8-141

(Yvon & al., 2004)

$$[x \setminus y \circ z \supset t] \iff t \in (y \bullet z) \setminus x$$

- **shuffle** $a \bullet b$ read sequences in a and b from left to right, allowing to switch from one string to another

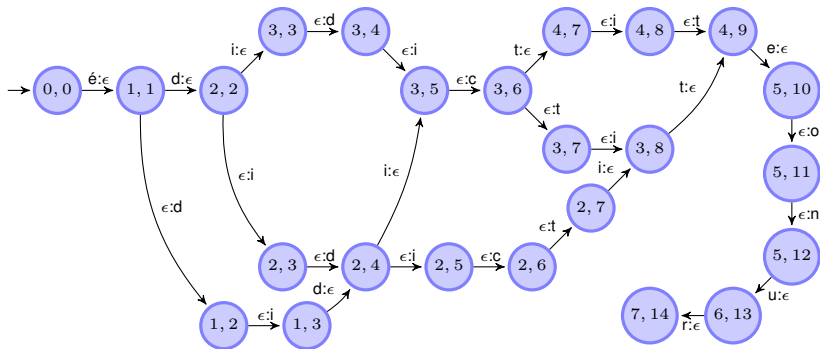
spondyodontilalgias et ondspondonylaltitisgia
 \in *spondylalgia* \circ *odontitis*

- **complement** $a \setminus b$ strings obtained by removing substring b in a

spondylitis \in *spondyodontilalgias* \setminus *odontalgia*
spydoniltis \in *spondyodontilalgias* \setminus *odontalgia*

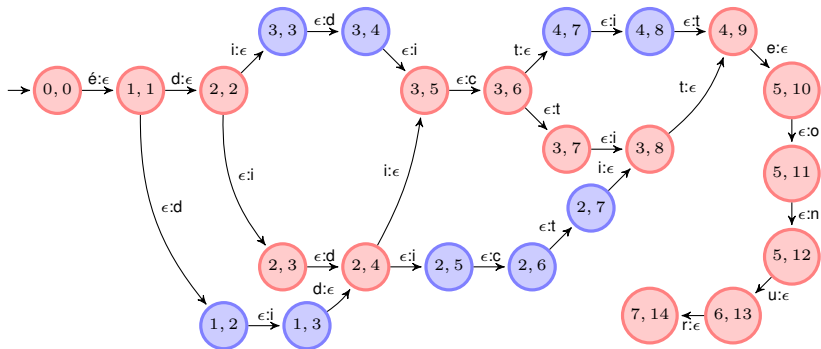
$\{(y \bullet z) \setminus x\}$ is a rational language

édition • dicteur \ éditeur



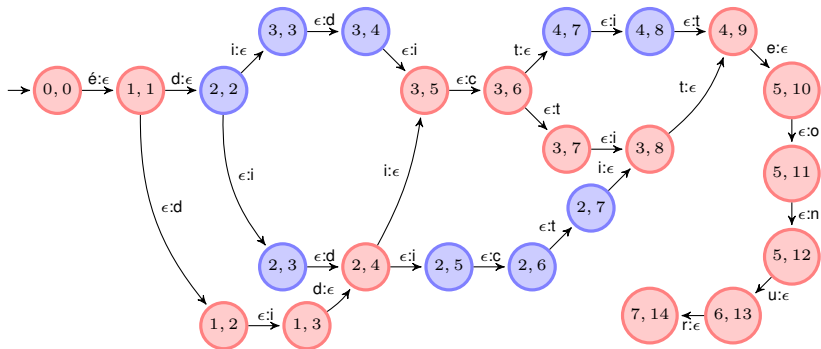
idction , diction , idcition, diiction, diciton

édition • dicteur \ éditeur



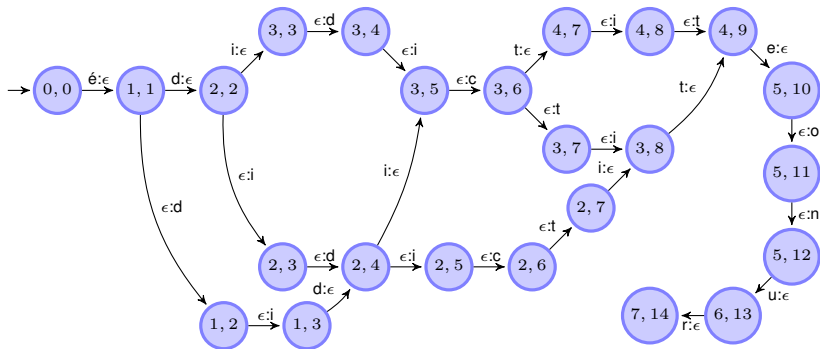
idction , diction , idcition, diiction, diciton

édition • dicteur \ éditeur



idction , **diction** , idciton, diicton, diciton

édition • dicteur \ éditeur



idction , diction , idciton, diicton, diciton

Solver alea [?]

- ▶ building the automaton typically is combinatorial
- ▶ alea randomly samples ρ shuffles of y and z and complements them (with x)
- ▶ **notes:**
 - ▶ for too low values of ρ , alea may fail to identify a solution
 - ▶ different shuffles may lead (after complementation) to the same solution
 - ▶ solutions are ranked in decreasing order of frequency

[*this guy drinks too much* \ *this boat sinks* \circ *those guys drink too much* \supset ?]

$\rho = 20$	$n = 8$
$t = 0.0003$	$rank = \phi$
<hr/>	
<i>thos_boate_sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>thoboatse__sinks</i> (2)	
<hr/>	
$\rho = 1000$	$n = 28$
$t = 0.009$	$rank = 2$
<hr/>	
<i>those_boat_ssink</i> (5)	
<i>those_boats_sink</i> (5)	
<i>thoes_tboa_sinks</i> (5)	
<hr/>	

$\rho = 100$	$n = 28$
$t = 0.001$	$rank = 13$
<hr/>	
<i>thoboatse__sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>those_sboat_sink</i> (2)	
<hr/>	
$\rho = 10^6$	$n = 19\,796$
$t = 3.82$	$rank = 10$
<hr/>	
<i>thoes_boat_sinks</i> (2550)	
<i>thoses_boat_sink</i> (1037)	
<i>those_boat_ssink</i> (999)	
<hr/>	

- ▶ ρ : nb of shuffles sampled
- ▶ n : # of solutions proposed

- ▶ $rank$: rank of the reference solution ($\phi \equiv$ not found)
- ▶ t : time (in sec.)

[*this guy drinks too much* \ *this boat sinks* \circ *those guys drink too much* \supset ?]

$\rho = 20$	$n = 8$
$t = 0.0003$	$rank = \phi$
<i>thos_boate_sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>thoboatse__sinks</i> (2)	

$\rho = 100$	$n = 28$
$t = 0.001$	$rank = 13$
<i>thoboatse__sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>those_sboat_sink</i> (2)	

$\rho = 1000$	$n = 28$
$t = 0.009$	$rank = 2$
<i>those_boat_ssink</i> (5)	
<i>those_boats_sink</i> (5)	
<i>thoes_tboa_sinks</i> (5)	

$\rho = 10^6$	$n = 19\,796$
$t = 3.82$	$rank = 10$
<i>thoes_boat_sinks</i> (2550)	
<i>thoses_boat_sink</i> (1037)	
<i>those_boat_ssink</i> (999)	

- ▶ ρ : nb of shuffles sampled
- ▶ n : # of solutions proposed

- ▶ $rank$: rank of the reference solution ($\phi \equiv$ not found)
- ▶ t : time (in sec.)

[*this guy drinks too much* \ *this boat sinks* \circ *those guys drink too much* \supset ?]

$\rho = 20$	$n = 8$
$t = 0.0003$	$rank = \phi$
<i>thos_boate_sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>thoboatse__sinks</i> (2)	

$\rho = 100$	$n = 28$
$t = 0.001$	$rank = 13$
<i>thoboatse__sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>those_sboat_sink</i> (2)	

$\rho = 1000$	$n = 28$
$t = 0.009$	$rank = 2$
<i>those_boat_ssink</i> (5)	
<i>those_boats_sink</i> (5)	
<i>thoes_tboa_sinks</i> (5)	

$\rho = 10^6$	$n = 19\,796$
$t = 3.82$	$rank = 10$
<i>thoes_boat_sinks</i> (2550)	
<i>thoses_boat_sink</i> (1037)	
<i>those_boat_ssink</i> (999)	

- ▶ ρ : nb of shuffles sampled
- ▶ n : # of solutions proposed

- ▶ $rank$: rank of the reference solution ($\phi \equiv$ not found)
- ▶ t : time (in sec.)

[*this guy drinks too much* \ *this boat sinks* ○ *those guys drink too much* ⊃ ?]

$\rho = 20$	$n = 8$
$t = 0.0003$	$rank = \phi$
<i>thos_boate_sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>thoboatse__sinks</i> (2)	

$\rho = 100$	$n = 28$
$t = 0.001$	$rank = 13$
<i>thoboatse__sinks</i> (2)	
<i>tho_boatse_sinks</i> (2)	
<i>those_sboat_sink</i> (2)	

$\rho = 1000$	$n = 28$
$t = 0.009$	$rank = 2$
<i>those_boat_ssink</i> (5)	
<i>those_boats_sink</i> (5)	
<i>thoes_tboa_sinks</i> (5)	

$\rho = 10^6$	$n = 19\,796$
$t = 3.82$	$rank = 10$
<i>thoes_boat_sinks</i> (2550)	
<i>thoses_boat_sink</i> (1037)	
<i>those_boat_ssink</i> (999)	

- ▶ ρ : nb of shuffles sampled
- ▶ n : # of solutions proposed

- ▶ $rank$: rank of the reference solution ($\phi \equiv$ not found)
- ▶ t : time (in sec.)

Solver word2vec

▶ $queen - king \sim woman - man$

[?]

▶ Solve $[x : y :: z : ?]$ by:

$$\hat{t} = \operatorname{argmax}_{t \in V} \cos(t, z - x + y)$$

▶ Can also solve **syntactic** equations:

$[work \setminus worked \circ accept \supset ?]$ ▶ accepted

▶ **Note:** rank words in V but does not generate new ones

Learning to Solve Equations

- ▶ making use of a training set of equations and their solutions
 $\mathcal{T} = \{((x, y, z), t)\}$ where $[x \setminus y \circ z \supset t]$

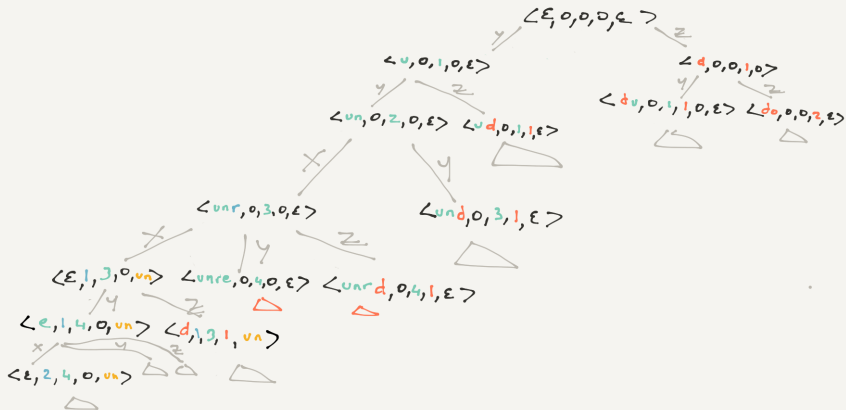
- ▶ using **structured learning**

- ▶ given $g : \mathcal{I}^3 \times \mathcal{I} \rightarrow \mathbb{R}$ which evaluates a fit between a triple of strings in \mathcal{I} , $i \equiv (x, y, z)$, and any string t , we seek to find (**search**):

$$\hat{t} = \operatorname{argmax}_{t \in \mathcal{I}} g(i, t)$$

- ▶ we assume a linear model for $g = \langle \mathbf{w}, \Phi(i, t) \rangle$ parametrized by a feature vector \mathbf{w} in \mathbb{R}^K and a **feature map** $\Phi(i, t)$ decomposed into K binary feature functions $\phi_k : (i, t) \rightarrow \{0, 1\}$.
- ▶ we wish to adjust \mathbf{w} so as to minimize over \mathcal{T} the number of search errors, thanks to the **voted perceptron** algorithm **[?]**

Search



Peculiarities of the Search

- ▶ Tree of depth $|x| + |y| + |z|$ with a branching factor close to 2 on average (too many nodes)
- ▶ Only few actions increase the prefix of a solution (in particular, many states with an empty prefix)

In practice, we:

- ▶ search space organized as a graph
- ▶ check that complementation is still possible
- ▶ control the maximum number of X or Y actions that can take place without complementing with x

Feature Map

3 families of **binary** features for characterizing $\langle s, i, j, k, p \rangle$:
 language model (14 features) evaluating the likelihood of the
 prefix p (or shuffle s) so far, according to a n-gram
 LM trained on an out-domain monolingual corpus.

- ▶ $\exists ?i : p_{LM}(p_i | p_{i-2} p_{i-1}) < \delta$

edit-distance (20 features) to enforce that the solution shares
 with y and z subsequences

- ▶ $[reader : unreadable :: doer : undoable]$

search-based (20k features)

- ▶ global: e.g. number of Y or Z actions in a row
- ▶ local: e.g. $(x_i, y_j, z_k) \equiv (r, a, d)$

Averaged Voted Perceptron [?]

$\mathbf{w}, \mathbf{w}_a \leftarrow \mathbf{0}$

$e \leftarrow 0$

repeat

$e \leftarrow e + 1$

for all example $(i, t) \in D$ **do**

$\hat{t} = \operatorname{argmax}_t \mathbf{w}^T \Phi(i, t)$

if $\hat{t} \neq t$ **then**

$\mathbf{w} \leftarrow \mathbf{w} + \Phi(i, t) - \Phi(i, \hat{t})$

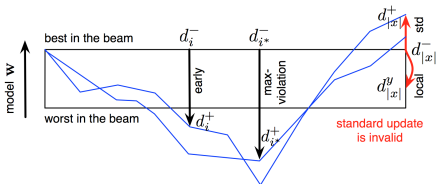
$\mathbf{w}_a \leftarrow \mathbf{w}_a + \mathbf{w}$

until converged

return $\mathbf{w}_a / e \cdot |\mathcal{D}|$

Averaged Voted Perceptron Variants

- ▶ Not suited for inexact search (our case) [?]
- ▶ update **only if** $\Phi(i, t) - \Phi(i, \hat{t}) < 0$
- ▶ We implemented 4 different variants (using forced-decoding)



standard do not care

safe remove the example if condition is not satisfied

early update on a prefix whenever there is no prefix of a reference solution in the beam

late update on the largest prefix that satisfies the condition



Metrics

accuracy percentage of test equations for which the solver output the correct solution at rank 1

silence percentage of test equations without any solution proposed

Dataset google

- Formal analogies of the dataset released by [?]

of analogies: 4977

word's avr. length: 7

adjective-adverbe	ADJ-ADV	[<i>amazing</i> \ <i>amazingly</i> ○ <i>serious</i> ⊃ <i>seriously</i>]
opposite	OPP	[<i>certain</i> \ <i>uncertain</i> ○ <i>competitive</i> ⊃ <i>uncompetitive</i>]
comparative	COMP	[<i>fast</i> \ <i>faster</i> ○ <i>bright</i> ⊃ <i>brighter</i>]
superlative	SUP	[<i>warm</i> \ <i>warmest</i> ○ <i>strange</i> ⊃ <i>strangest</i>]
present-participle	PP	[<i>code</i> \ <i>coding</i> ○ <i>dance</i> ⊃ <i>dancing</i>]
nationality-adverb	NAT	[<i>Australia</i> \ <i>Australian</i> ○ <i>Croatia</i> ⊃ <i>Croatian</i>]
past-tense	PAST	[<i>decreasing</i> \ <i>decreased</i> ○ <i>listening</i> ⊃ <i>listened</i>]
plural	PLUR	[<i>eye</i> \ <i>eyes</i> ○ <i>donkey</i> ⊃ <i>donkeys</i>]
plural-verbs	PL-VB	[<i>listen</i> \ <i>listens</i> ○ <i>eat</i> ⊃ <i>eats</i>]

Dataset: m s r

- Formal analogies of the dataset released by [?]

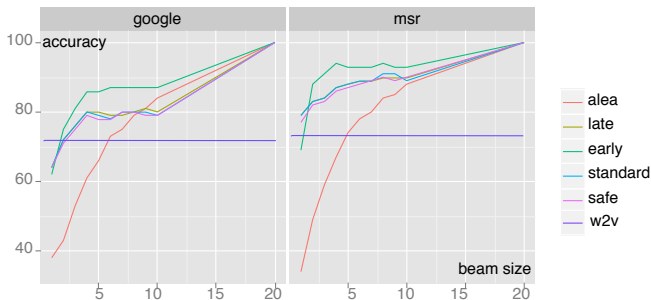
of analogies: 3664

word's avr. length: 6

JJ-JJR	[<i>high \ higher</i> ○ <i>wild</i> ⊃ <i>wilder</i>]
JJR-JJ	[<i>greater \ greatest</i> ○ <i>earlier</i> ⊃ <i>earliest</i>]
JJS-JJ	[<i>low \ lowest</i> ○ <i>short</i> ⊃ <i>shortest</i>]
NN-NNPOS	[<i>problem \ problems</i> ○ <i>program</i> ⊃ <i>programs</i>]
VB-VBP	[<i>take \ takes</i> ○ <i>run</i> ⊃ <i>runs</i>]
VB-VBD	[<i>prevent \ prevented</i> ○ <i>consider</i> ⊃ <i>considered</i>]
NNPOS-NN	[<i>days \ day</i> ○ <i>citizens</i> ⊃ <i>citizen</i>]
VBZ-VBD	[<i>believes \ believed</i> ○ <i>likes</i> ⊃ <i>liked</i>]

Accuracy as a function of the beam size (or ρ)

- ▶ training on `msr` and testing on `google` (left), or the other way round (right)

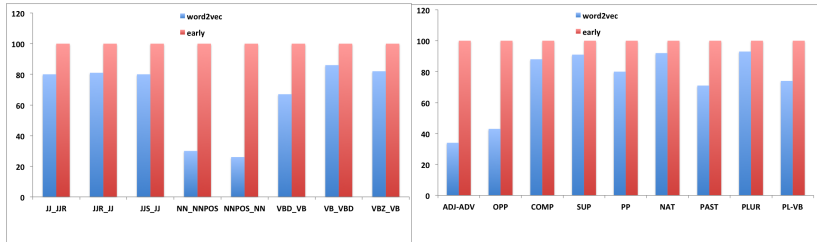


A Comparison to State-of-the-art

- trust `early` if the solution defines with the equation a formal analogy, `word2vec` otherwise

	<code>msr</code>	<code>google</code>
<code>word2vec</code>	67%	63%
<code>early+word2vec</code>	72%	71%
Levy et al., 2015	72.9%	75.8%

Comparison by Types



Dataset

► Much harder to get. We did this:

1 Train a phrasal translation table (from an EN-ES parallel corpus)

a actualizar los acuerdos | to update the agreements | 0.0047
a cambiar la base | to change the basis | 0.0035545
basado en el trabajo de | based on the efforts of | 2.02579e-

2 Split good enough associations into R a reference set, and M a translation memory

3 For each $(s, e) \in R$,

- use an analogical device to translate s into e , using M as a memory, that is, finding $(x, x'), (y, y'), (z, z') \in M$ such that:

$$[x \setminus y \circ z \supset s] \text{ and } [x' \setminus y' \circ z' \supset e]$$

- collecting $[x' \setminus y' \circ z' \supset e]$ (analogies in the English side)

4 Select:

- simple: a random subset of those equations
- hard: a sampling of equations more difficult for `alea` to solve

Dataset: simple

simple

phrase's avr. length: 16

[<i>international investigation</i>	\	<i>international democracy</i>	○
	<i>an international investigation</i>	▷	<i>an international democracy</i>]
[<i>young girls</i>	\	<i>training of young girls</i>	○
	<i>young girls and</i>	▷	<i>training of young girls and</i>]
[<i>political situation is viable</i>	\	<i>political situation is still</i>	○
	<i>the political situation is viable</i>	▷	<i>the political situation is still</i>]

Dataset: hard

hard

phrase's avr. length: 17

- | | | | | |
|---|------------------------------------|---|------------------------------------|---|
| [| <i>adopted recently by</i> | \ | <i>recently adopted by</i> | ○ |
| | <i>study published recently by</i> | ▷ | <i>study recently published by</i> |] |
| [| <i>competition and the</i> | \ | <i>competition and against</i> | ○ |
| | <i>competition and of the</i> | ▷ | <i>of competition and against</i> |] |
| [| <i>their governments to</i> | \ | <i>their governments are</i> | ○ |
| | <i>their governments and to</i> | ▷ | <i>and their governments are</i> |] |

Accuracy (Silence)

`test` \equiv `simple` Structured solvers trained on 10 epochs with $\eta = 7$, against `alea` with $\rho = 1000$

train	simple	hard
standard	30.6 (7)	30.1 (10)
early	38.4 (8)	27.6 (10)
late	26.9 (9)	20.3 (11)
safe	25.4 (8)	25.6 (13)
alea	33 (0)	

- ▶ `alea` is competitive, but outperformed by `early`
- ▶ preferable to train on simple analogies

Accuracy (Silence)

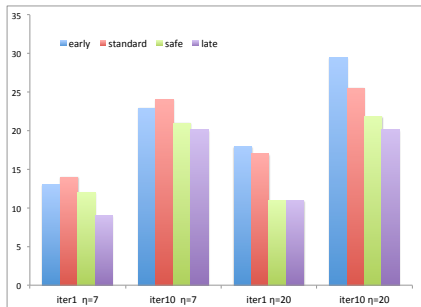
`test` \equiv `hard` Structured solvers trained on 10 epochs with
 $\eta = 7$, against `alea` with $\rho = 1000$

<code>train</code>	<code>simple</code>	<code>hard</code>
<code>standard</code>	19.6 (17)	24.0 (56)
<code>early</code>	26.0 (6)	22.9 (7)
<code>late</code>	18.9 (9)	20.1 (6.7)
<code>safe</code>	14.6 (18.7)	21.0 (56)
<code>alea</code>	18 (0)	

- ▶ more challenging dataset !
- ▶ structured learning systematically better
- ▶ preferable to train on harder to solve analogies
- ▶ silence rate rather high

Accuracy on `hard`

- ▶ after 1 and 10 epochs
- ▶ for 2 values of η : 7 and 20



- ▶ more epochs is (of course) preferable
- ▶ opening the search space as well (expectedly)
- ▶ `early` seems overall the best variant tested

The background features four circular segments, each composed of ten colored triangles pointing towards the center. The colors used are blue, green, yellow, red, orange, and purple. The segments are arranged in a 2x2 grid pattern around the central text.

Conclusion

Futur

En cours

- ▶ probabiliser le solveur (thèse de Rafik Rhouma)

Long terme

- ▶ probabiliser tout le processus (search, solveur, agrégation)
- ▶ passage à l'échelle
- ▶ noyau analogique ?
- ▶ analogie sur les arbres [?, ?]

Bibliography I