

# Candidate Evaluation Strategies for Improved Difficulty Prediction of Language Proficiency Tests



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



**DIPF**

Educational Research  
and Educational Information

**BEA Workshop 2015**

**Lisa Beinborn, Torsten Zesch, Iryna Gurevych**

**UKP Lab**

**Department of Computer Science**

**Technische Universität Darmstadt**



# Follow-up work

---

TACL paper 2014:

## “Predicting the Difficulty of Language Proficiency Tests”

What is difficult for language learners? (And why?)

How can we predict difficulty?

Data: Language Proficiency Tests

Task: Predict difficulty of the items

# Reduced Redundancy Principle

- Spolsky (1969): “Natural language is redundant”
- The ability to deal with reduced redundancy distinguishes beginners from advanced language learners

Thanks to the redundancy of language,  
yxx cxn xndxrstxnd whxt x xm wrxtxng  
xvxn xf x rxplxcx xll thx vxwxls wxth xn ‘x’.

t gts lttl hrdr f y dn't  
vn kn whr th vwls r.

Steven Pinker, *The Language Instinct: How the Mind Creates Language* (William Morrow, 1994)

## C-Test [Klein-Braley and Raatz, 1982]



- Beginning and end of text provide context
- Every second word is a gap
- Smaller “half” of the word is provided

The roots of humanity can be traced back to millions of years ago. T\_\_\_\_\_ primary evid\_\_\_\_\_ comes fr\_\_\_\_\_ fossils - skulls, skel\_\_\_\_\_ and bo\_\_\_\_\_ fragments. Scien\_\_\_\_\_ have ma\_\_\_\_\_ tools th\_\_\_\_\_ allow th\_\_\_\_\_ to ext\_\_\_\_\_ subtle infor\_\_\_\_\_ from anc\_\_\_\_\_ bones a\_\_\_\_\_ their enviro\_\_\_\_\_ settings. Mod\_\_\_\_\_ forensic wo\_\_\_\_\_ in t\_\_\_\_\_ field a\_\_\_\_\_ in labora\_\_\_\_\_ can n\_\_\_\_\_ provide a rich understanding of how our ancestors lived.

## C-Test [Klein-Braley and Raatz, 1982]



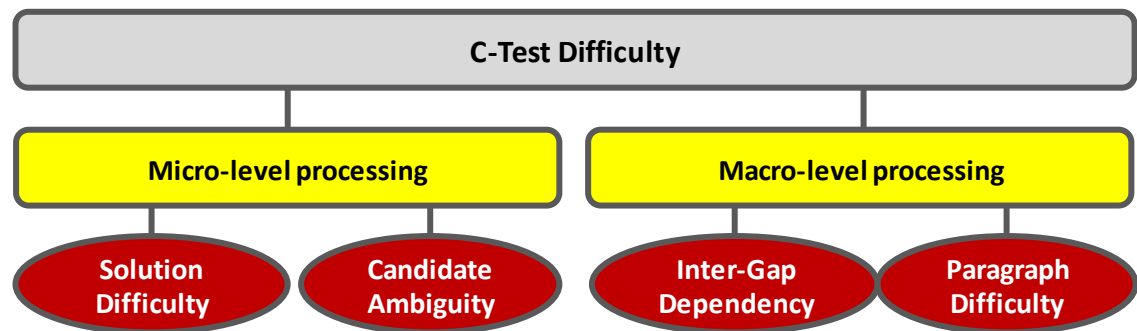
- Beginning and end of text provide context
- Every second word is a gap
- Smaller “half” of the word is provided

The roots of humanity can be traced back to millions of years ago. T\_\_\_\_\_ primary evid\_\_\_\_\_ comes fr\_\_\_\_\_ fossils - skulls, skel\_\_\_\_\_ and bo\_\_\_\_\_ fragments. Scien\_\_\_\_\_ have ma\_\_\_\_\_ tools th\_\_\_\_\_ allow th\_\_\_\_\_ to ext\_\_\_\_\_ subtle infor\_\_\_\_\_ from anc\_\_\_\_\_ bones a\_\_\_\_\_ their enviro\_\_\_\_\_ settings. Mod\_\_\_\_\_ forensic wo\_\_\_\_\_ in t\_\_\_\_\_ field a\_\_\_\_\_ in labora\_\_\_\_\_ can n\_\_\_\_\_ provide a rich understanding of how our ancestors lived.

**Data: TU Darmstadt  
Placement test at  
language centre**

# Yesterday: Difficulty Prediction of English C-test Gaps

*The roots of humanity can be traced back to millions of years ago. The primary evidence comes from fossils - skulls, skeletal and bone fragments. Scientists have made tools that allow them to extract subtle information from ancient bones and their environmental settings. Modern forensic work in the field and in the laboratory can now provide a rich understanding of how our ancestors lived.*



# Outlook of TACL paper

- 1) Adapt to other languages and to other test variants
- 2) Improve the dimensions candidate ambiguity and inter-gap dependency

# TODAY: Part 1



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

- 1) Adapt to **German and French** and to **two test variants**
- 2) Improve the dimensions candidate ambiguity and inter-gap dependency



# TODAY, Part 2



- 1) Adapt to German and French and to two **test variants**
- 2) Improve the dimensions **candidate ambiguity** and inter-gap dependency

## Was ist Kreativität?

Die meisten halten Kreativität für eine seltene Gabe, über die nur eine exklusive Minderheit verfügt. Dabei \_\_\_\_\_t jeder \_\_\_\_\_sch auf \_\_\_\_\_ne Weise \_\_\_\_\_erisch. Wir \_\_\_\_\_nen nicht \_\_\_\_\_cht kreativ \_\_\_\_\_in. Die \_\_\_\_\_idende Frage \_\_\_\_\_t, ob \_\_\_\_\_r diese \_\_\_\_\_liche Fähigkeit \_\_\_\_\_iv pflegen \_\_\_\_\_er verkümmern \_\_\_\_\_sen. Denn \_\_\_\_\_vität bezieht \_\_\_\_\_ch nicht \_\_\_\_\_f ein \_\_\_\_\_mmtes Themengebiet, \_\_\_\_\_em ist \_\_\_\_\_all möglich. \_\_\_\_\_ativ sein \_\_\_\_\_utet, sich \_\_\_\_\_as anderes \_\_\_\_\_ellen zu \_\_\_\_\_nen als das, was man gerade sieht. Kreativität wird häufig mit Innovation verwechselt.

13. His characteristic talk , with its keen \_\_\_\_\_ of detail and subtle power of inference held me amused and enthralled. :

- instincts
- presumption
- observance
- expiation
- implements

# Different Languages

Is it a language-independent approach?

Can it be adapted to other languages?



Markus Koljonen: [iki.fi/markus.koljonen](http://iki.fi/markus.koljonen)



## Le Noël des familles

Noël semble immuable, et pourtant il change ! Au co\_\_\_\_ du te\_\_\_\_ on s'aper\_\_\_\_, discrètement ma\_\_\_\_ sûrement, q\_\_\_\_ la fê\_\_\_\_ de l\_\_\_\_ Nativité pr\_\_\_\_ de nouv\_\_\_\_ allures e\_\_\_\_ alliant l\_\_\_\_ traditions e\_\_\_\_ la mode\_\_\_\_. Traditionnellement, l\_\_\_\_ sapin d\_\_\_\_ Noël e\_\_\_\_ surtout déc\_\_\_\_ de bou\_\_\_\_, de guirl\_\_\_\_ lumineuses e\_\_\_\_ de boules multicolores. Toutefois, une certaine influence germanique se fait de plus en plus sentir, avec la présence de figurines en bois ou de dessins aux fenêtres.

- articles
- accents
- richer morphology



## Le Noël des familles

Noël semble immuable, et pourtant il change ! Au co\_\_\_\_ du te\_\_\_\_ on s'aper\_\_\_\_, discrètement ma\_\_\_\_ sûrement, q\_\_\_\_ la fê\_\_\_\_ de l\_\_\_\_ Nativité pr\_\_\_\_ de nouv\_\_\_\_ allures e\_\_\_\_ alliant l\_\_\_\_ traditions e\_\_\_\_ la mode\_\_\_\_. Traditionnellement, l\_\_\_\_ sapin d\_\_\_\_ Noël e\_\_\_\_ surtout déc\_\_\_\_ de bou\_\_\_\_, de guirl\_\_\_\_ lumineuses e\_\_\_\_ de boules multicolores. Toutefois, une certaine influence germanique se fait de plus en plus sentir, avec la présence de figurines en bois ou de dessins aux fenêtres.

**Data: TU Darmstadt  
Placement test at  
language centre**



Auf einer Weltkarte kann man sehen, dass Asien und Nordamerika im Norden nur durch einen schmalen Meeresstreifen voneinander getrennt sind: durch die Beringstraße. Während d\_\_\_\_\_ Eiszeit herrs\_\_\_\_\_ auf d\_\_\_\_\_ ganzen Er\_\_\_\_\_ niedrigere Temper\_\_\_\_\_, und d\_\_\_\_\_ Beringstraße w\_\_\_\_\_ zugefroren. Üb\_\_\_\_\_ das E\_\_\_\_\_ gelangten Volksg\_\_\_\_\_ aus As\_\_\_\_\_ auf d\_\_\_\_\_ amerikanischen Kont\_\_\_\_\_. Manche bli\_\_\_\_\_ in Norda\_\_\_\_\_ und bild\_\_\_\_\_ mehr a\_\_\_\_\_ tausend versch\_\_\_\_\_ Stämme - jew\_\_\_\_\_ mit ei\_\_\_\_\_ eigenen Sprache, die anderen zogen weiter bis nach Südamerika.

- articles, cases
- Umlaute, case sensitivity
- old/new orthography
- compounds

# German C-Test



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Auf einer Weltkarte kann man sehen, dass Asien und Nordamerika im Norden nur durch einen schmalen Meeresstreifen voneinander getrennt sind: durch die Beringstraße. Während d\_\_\_\_\_ Eiszeit herrs\_\_\_\_\_ auf d\_\_\_\_\_ ganzen Er\_\_\_\_\_ niedrigere Temper\_\_\_\_\_, und d\_\_\_\_\_ Beringstraße w\_\_\_\_\_ zugefroren. Üb\_\_\_\_\_ das E\_\_\_\_\_ gelangten Volksg\_\_\_\_\_ aus As\_\_\_\_\_ auf d\_\_\_\_\_ amerikanischen Kont\_\_\_\_\_. Manche bli\_\_\_\_\_ in Norda\_\_\_\_\_ und bild\_\_\_\_\_ mehr a\_\_\_\_\_ tausend versch\_\_\_\_\_ Stämme - jew\_\_\_\_\_ mit ei\_\_\_\_\_ eigenen Sprache, die anderen zogen weiter bis nach Südamerika.

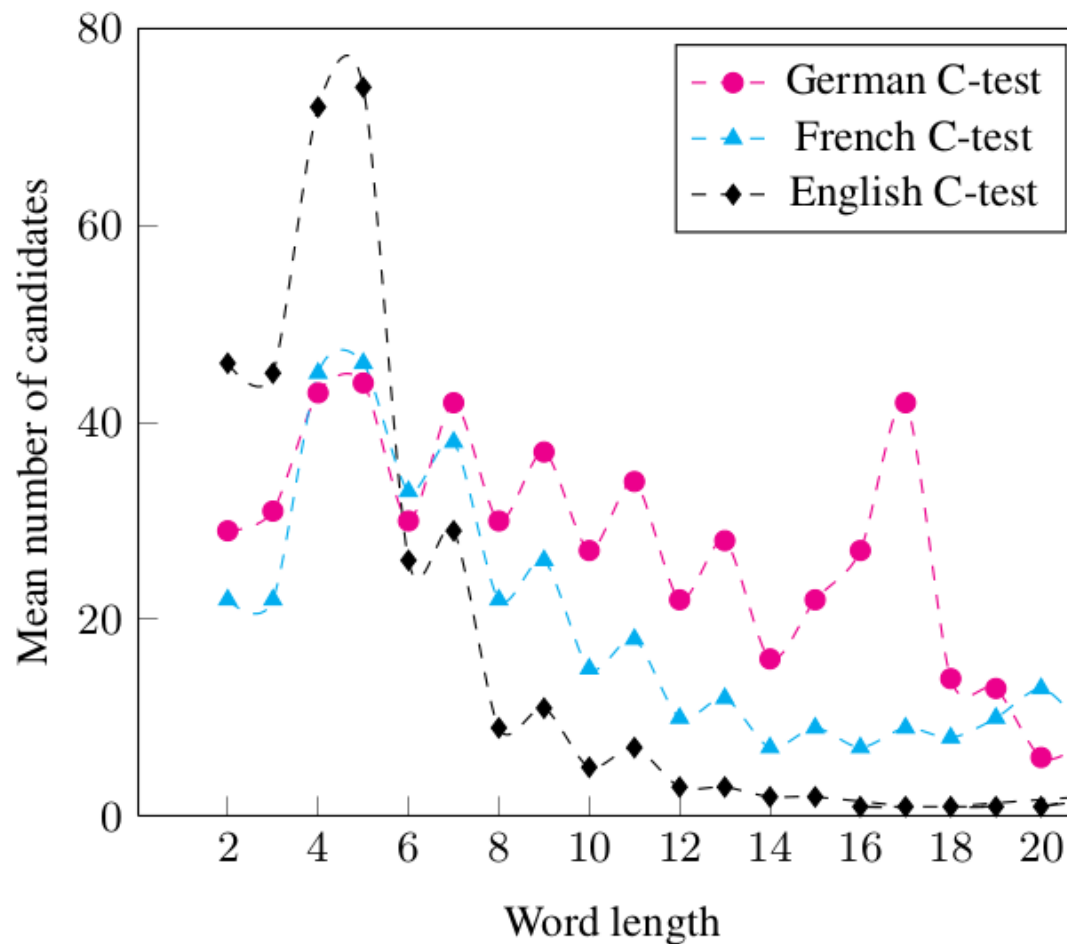
**Data: TestDaF Institute  
Certificate of German  
proficiency for  
university admission**

# Different Languages...



Language	Words	Mean word length
English	99,171	8.5 ± 2.6
French	139,719	9.6 ± 2.6
German	332,263	12.0 ± 3.5

# Different Candidate Space





# Different Test Types

---



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

**The reduced redundancy principle comprises more than just the C-test.**

# Prefix Deletion Test



Was ist Kreativität?

Die meisten halten Kreativität für eine seltene Gabe, über die nur eine exklusive Minderheit verfügt. Dabei \_\_\_\_\_t jeder \_\_\_\_\_sch auf \_\_\_\_\_ne Weise \_\_\_\_\_erisch. Wir \_\_\_\_\_nen nicht \_\_\_\_\_cht kreativ \_\_\_\_\_in. Die \_\_\_\_\_idende Frage \_\_\_\_\_t, ob \_\_\_\_\_r diese \_\_\_\_\_liche Fähigkeit \_\_\_\_\_iv pflegen \_\_\_\_\_er verkümmern \_\_\_\_\_sen. Denn \_\_\_\_\_vität bezieht \_\_\_\_\_ch nicht \_\_\_\_\_f ein \_\_\_\_\_mmtes Themengebiet, \_\_\_\_\_ern ist \_\_\_\_\_all möglich. \_\_\_\_\_ativ sein \_\_\_\_\_utet, sich \_\_\_\_\_as anderes \_\_\_\_\_ellen zu \_\_\_\_\_nen als das, was man gerade sieht. Kreativität wird häufig mit Innovation verwechselt.

- **prefix vs postfix**
- **multiple solutions**

# Prefix Deletion Test



Was ist Kreativität?

Die meisten halten Kreativität für eine seltene Gabe, über die nur eine exklusive Minderheit verfügt. Dabei \_\_\_\_\_t jeder \_\_\_\_\_sch auf \_\_\_\_\_ne Weise \_\_\_\_\_erisch. Wir \_\_\_\_\_nen nicht \_\_\_\_\_cht kreativ \_\_\_\_\_in. Die \_\_\_\_\_idende Frage \_\_\_\_\_t, ob \_\_\_\_\_r diese \_\_\_\_\_liche Fähigkeit \_\_\_\_\_iv pflegen \_\_\_\_\_er verkümmern \_\_\_\_\_sen. Denn \_\_\_\_\_vität bezieht \_\_\_\_\_ch nicht \_\_\_\_\_f ein \_\_\_\_\_mmtes Themengebiet, \_\_\_\_\_ern ist \_\_\_\_\_all möglich. \_\_\_\_\_ativ sein \_\_\_\_\_utet, sich \_\_\_\_\_as anderes \_\_\_\_\_ellen zu \_\_\_\_\_nen als das, was man gerade sieht. Kreativität wird häufig mit Innovation verwechselt.

**Data: University of  
Duisburg-Essen  
German proficiency test for  
prospective teachers**



13. His characteristic talk , with its keen \_\_\_\_\_ of detail and subtle power of inference held me amused and enthralled.

- instincts
- presumption
- observance
- expiation
- implements

- Closed format
- Distractors

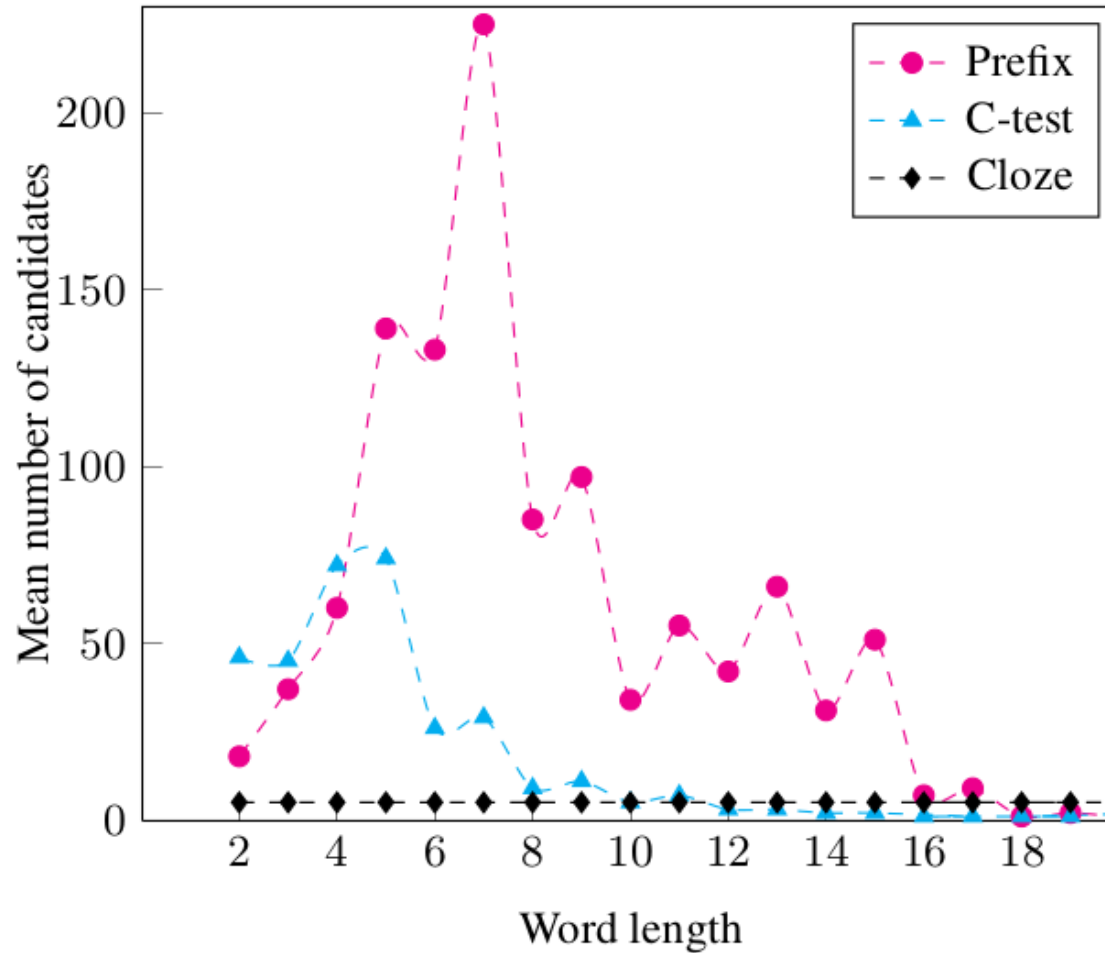


13. His characteristic talk , with its keen \_\_\_\_ of detail and subtle power of inference held me amused and enthralled.

- instincts
- presumption
- observance
- expiation
- implements

**Questions:**  
**Microsoft Research**  
**Zweig & Burges (2012)**  
**Error Rates: own study**

# Candidate Space



# Data



Format	Test Type	Texts	Gaps	Participants	Av. Error Rate
Open	C-test en	39	775	210	.35 ± .25
	C-test fr	40	799	24	.52 ± .28
	C-test de	82	1,640	251	.55 ± .26
	Prefix de	14	348	225	.36 ± .23
Closed	Cloze en	100	100	22	.27 ± .22

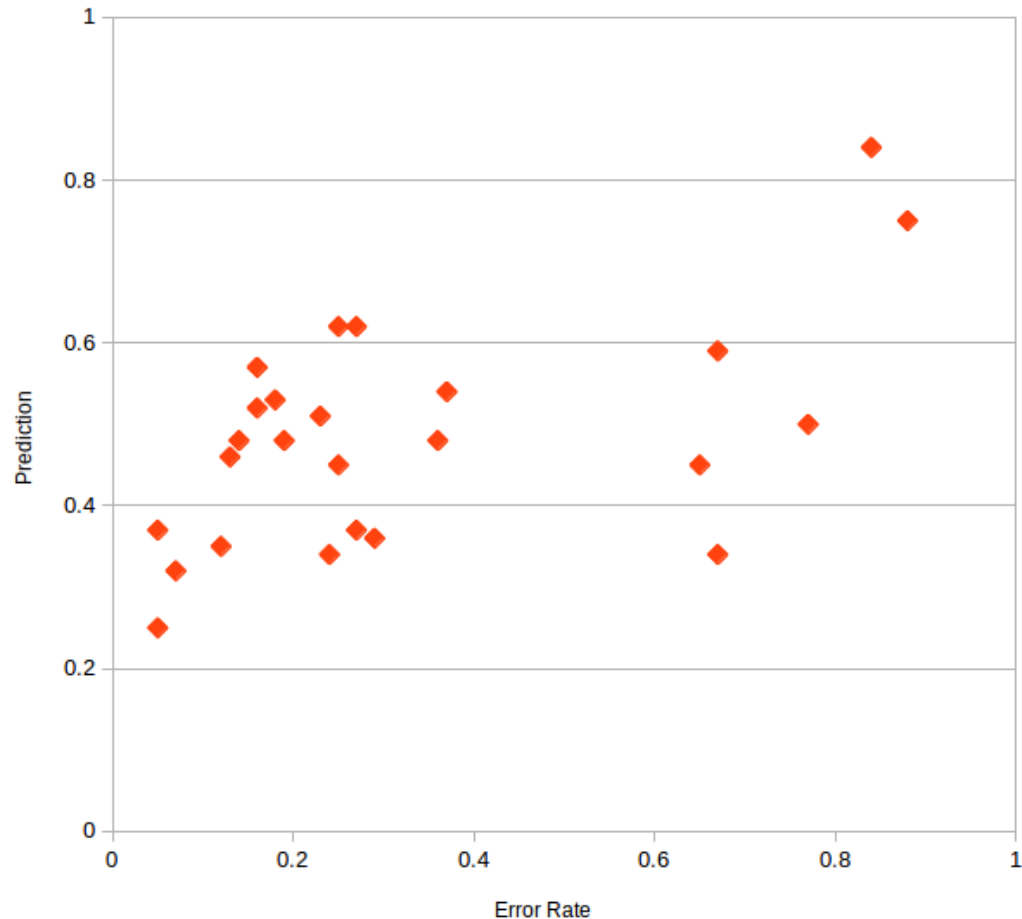
# Generalization

- Adapt format and types
- Adapt linguistic pre-processing (DKPro)
- Adapt resources
- Adapt features (Reduction from 87 to 70)



# Prediction Task

- SMO regression, Leave-one-out cross-validation on texts



# Prediction Results: Languages



Test Type	Pearson's r
C-test en	.47
C-test fr	.67
C-test de	.61

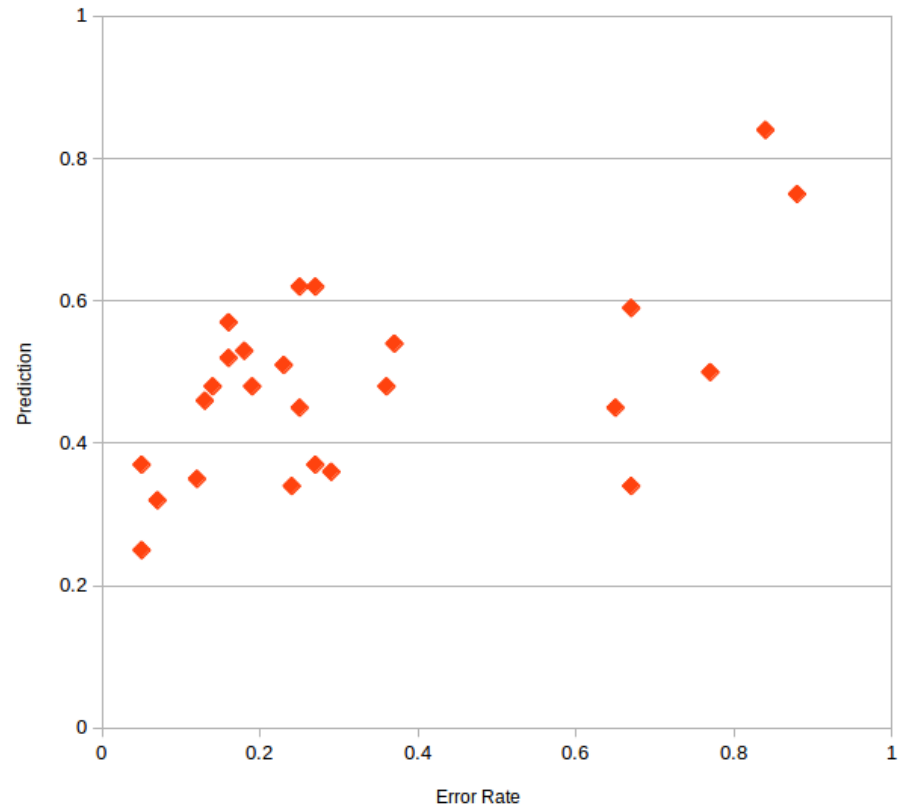
# Prediction Results



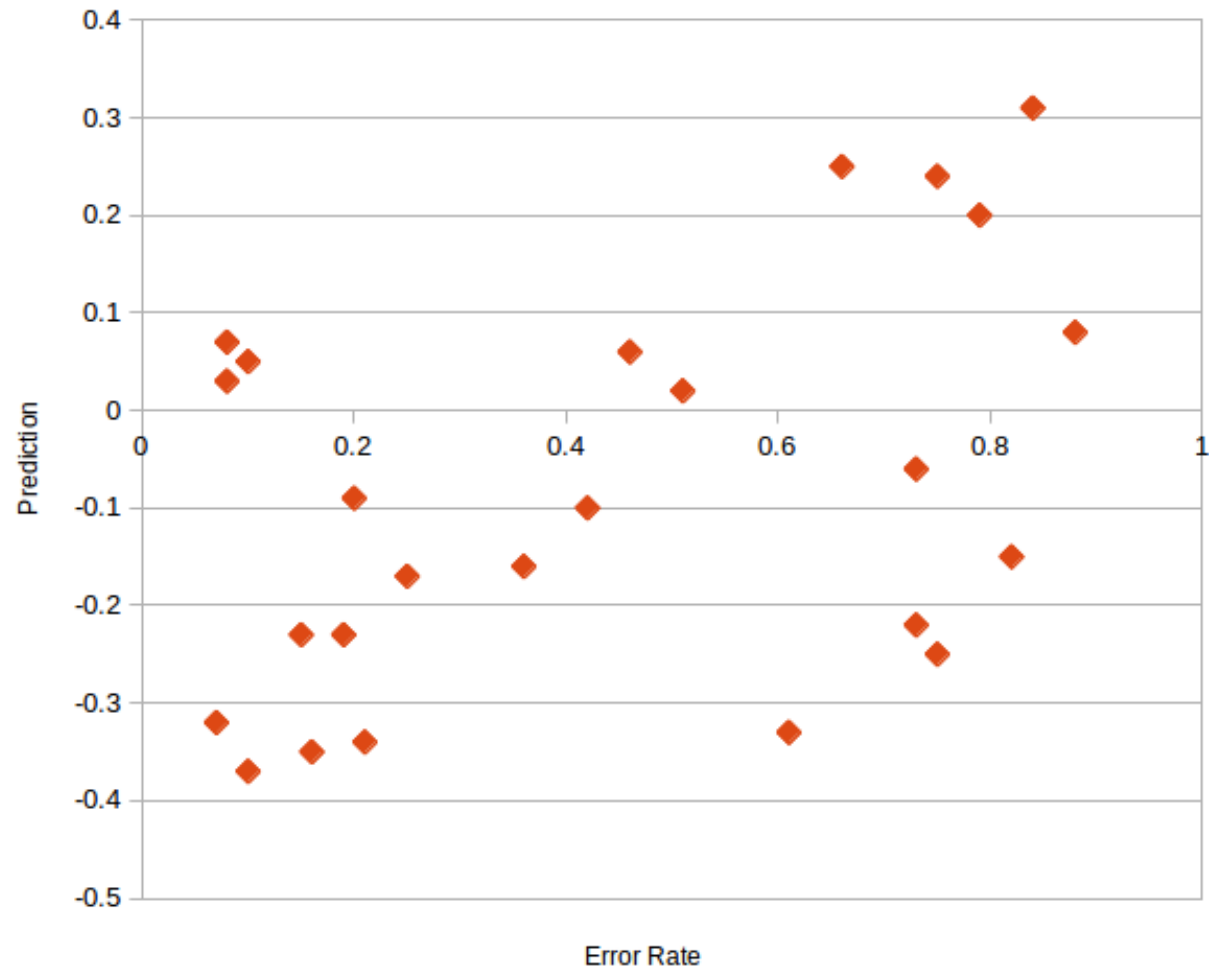
Test Type	Pearson's r
C-test en	.47
C-test fr	.67
C-test de	.61
Prefix de	.27
Cloze en	.20

# Prediction Results

Test Type	Pearson's r
C-test en	.47
C-test fr	.67
C-test de	.61
Prefix de	.27
Cloze en	.20



# Prefix Deletion: Outlier



# Open vs closed format

- Many of the features developed for difficulty prediction were targeted at production problems (e.g. spelling).
- Closed cloze tests only require recognition skills.

Test Type	Pearson's r
C-test en	.47
C-test fr	.67
C-test de	.61
Prefix de	.27
Cloze en	.20

# A closer look at cloze



49. The stage lost a fine \_\_\_\_\_, even as science lost an acute reasoner, when he became a specialist in crime.

- linguist
- hunter
- actor
- estate
- horseman

**Error Rate: 0.0**



21. When his body had been carried from the cellar we found ourselves still confronted with a problem which was almost as \_\_\_\_\_ as that with which we had started.

- tall
- loud
- invisible
- quick
- formidable

**Error Rate: 0.6**



# Inspiration from automated solving

How would we proceed for nlp-based solving of language tests?

Rank the candidates

- 1) According to language model probability
- 2) According to semantic relatedness between candidate and context



# Differences

- Automatic solving:  
train on domain-specific data (Holmes novels)
- Difficulty prediction:  
train on general data to model learner knowledge

Language model: trained with Berkeley LM on Leipzig corpora

- 1 million sentences, news domain

Explicit Semantic Analysis Index: trained on Wikipedia

- 931,559 concepts

# LM ranker

- Candidate fitness: log probability of sentence in language model

The stage lost a fine \_\_\_\_\_ , even as science lost an acute reasoner, when he became a specialist in crime .

<b>actor</b>	-358.83
estate	-361.22
hunter	-361.96
linguist	-362.71
horseman	-362.93

Rank of solution: **1**

- Candidate fitness: log probability of sentence in language model

When his body had been carried from the cellar we found ourselves still confronted with a problem which was almost as \_\_\_\_\_ as that with which we had started .

tall	-175.17
quick	-176.61
loud	-178.60
<b>formidable</b>	-179.52
invisible	-179.52

Rank of solution: 4

- Candidate fitness: sum over the cosine similarity between the candidate and every content word in the sentence (similar to Zweig et al. 2012)

The stage lost a fine \_\_\_\_\_ , even as science lost an acute reasoner, when he became a specialist in crime .

<b>actor</b>	1.06
estate	0.70
hunter	0.58
horseman	0.29
linguist	0.29

Rank of solution: **1**

# ESA ranker

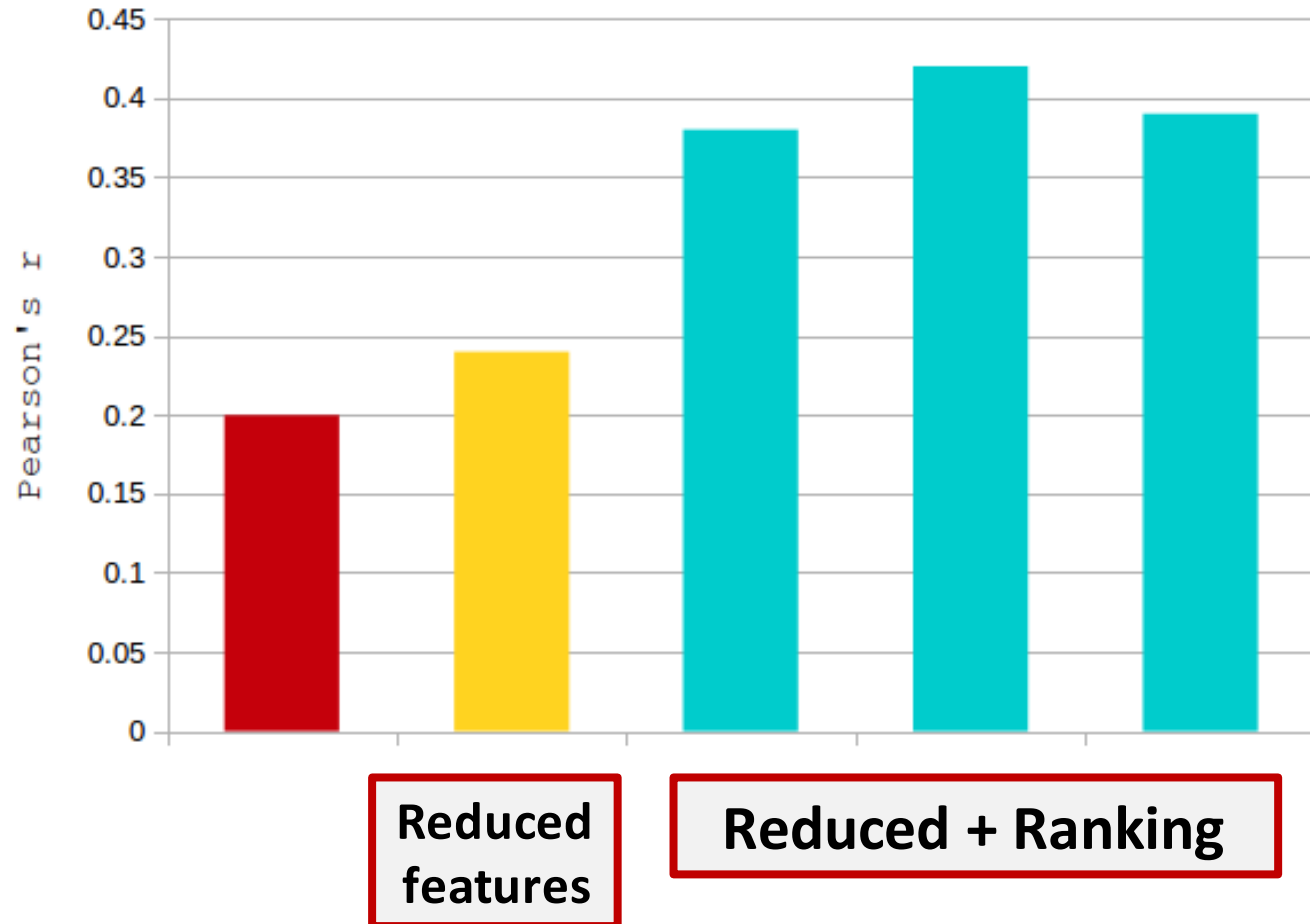
- Candidate fitness: sum over the cosine similarity between the candidate and every content word in the sentence (similar to Zweig et al. 2012)

When his body had been carried from the cellar we found ourselves still confronted with a problem which was almost as \_\_\_\_\_ as that with which we had started.

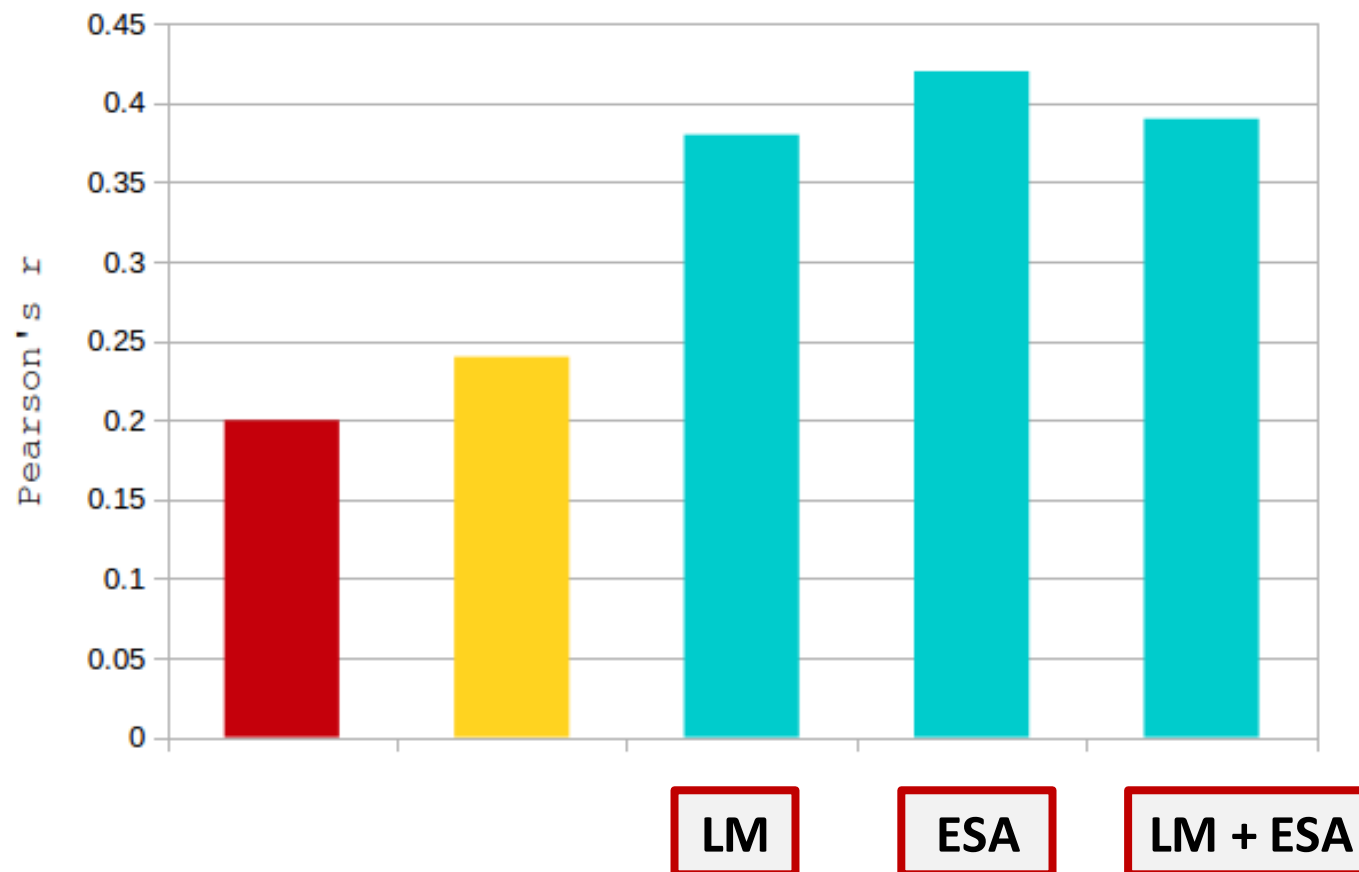
quick	1.03
tall	0.80
invisible	0.63
loud	0.61
<b>formidable</b>	0.56

Rank of solution: 5

# Improved Cloze Results



# Improved Cloze Results



# Summary and Outlook

- Successfully adapted the prediction framework to new languages and test types.
- We ❤️ French.
- Candidate evaluation strategies from automatic solving can be applied to simulate learner difficulties.

## Plans:

- Currently collecting more error rates for cloze tests.
- Use candidate evaluation strategies for distractor generation and difficulty manipulation.
- Strategies for inter-item dependencies.





**Thank you very much for your \_\_\_\_\_.**

- attention
- presence
- patience
- enthusiasm

**Do you have any \_\_\_\_\_?**

- questions
- comments
- praises
- doubts

