## Information Density and Scientific Literacy in English

Stefania Degaetano-Ortlieb, Hannah Kermes, Ashraf Khamis, Noam Ordan, Elke Teich



## Aim

Investigate

- in general: formation of discourse types/registers
- specific: diachronic development of written scientific English

### Motifs

- Specialization (greater encoding density)
- Conventionalization (greater linguistic uniformity)

Balance in Information Density (ID)

## Example: Abstract



basic ling structure: complex NPs, simple clause structure

- basic linguistic structure: complex NPs, simple clause structure
- local effects: omission of determiners
- global effects: high TTR, high lexical density

The Excesses of the Sines of the Refraction of several sorts of Rays above their common Sine of Incidence when the Refractions are made out of divers denser Mediums immediately into one and the same rarer Medium, suppose of Air, are to one another in a given Proportion. We report the discovery of a novel downstream target of BCR-ABL signalling, PRL-3 (PTP4A3), an oncogenic tyrosine phosphatase. Analysis of CML cancer cell lines and CML patient samples reveals the upregulation of PRL-3.

## Assumptions and Hypotheses

- Higher encoding density over time (cf. Halliday, 1989: On the language of physical science)
  - scientific texts will exhibit higher ID *relative* to other productions / "general language" over time
  - within scientific productions, ID will decrease over time
- $\rightarrow$  specialization
  - $\rightarrow$  conventionalization
- Correlation between variation in ling encoding and ID
  - Ling features marking specialization and conventionalization serve optimizing ID in scientific writing
    - longer, expanded ling forms  $\rightarrow$  less predictable, more informative
    - shorter, reduced ling forms  $\rightarrow$  more predictable, less informative

Corpus-based approach:Find out linguistic features indicative of IDLM approach:Measure ID

- synchronically by comparison of Abstracts vs. Research articles
   Corpus: SciTex, size: ~34M words, 9 academic disciplines, divisions annotation (two time periods 70/80s and 2000s)
- diachronically by comparison of historical stages of text productions (17th century - present)

Corpus: RSC, size: ~23M words, time periods: 1665-1870

Corpus-based approach:

- Feature selection, extraction, evaluation
- Classification with SVM

### Synchronic analysis

- Abstracts vs. Research articles
   Diachronic analysis
- Historical stages

higher encoding density	lower encoding density
high TTR	low TTR
high TTR on lexical words	low TTR on lexical words
high lexical density	low lexical density
complex NPs	simple NPs
many nominalizations	few nominalizations
few formulaic expressions	many formulaic expressions
simple sentence structure ('x is y')	more varied sentence structures
many nominal compounds	few nominal compounds
few determiners (omissions)	many determiners
few relativizers (omissions)	many relativizers
few complementizers (omissions)	many complementizers

LM-based approach:

- ID in single texts
  - Cross-entropy based
     on Genzel and Charniak (2002))
  - Calculate entropy at each token position
  - Explore idea of entropy rate constancy
  - Sliding window of 4 tokens
  - Simple model of memory decay

(by Peter Fankhauser, IDS Mannheim)

Most important for nuclear reaction studies are Van de Graaff accelerators in which ions are accelerated in an evacuated tube by an electrostatic real maintained between a high voltage terminal and an earth terminal charge being conveyed to a high voltage terminal by a rotating belt or chain as

In early forms a this accelerator, positive ions from a **Gaseous discharge** tube were accelerated from a high voltage terminal to earth at But, in modern **tandem** accelerators, **negative** ions at accelerated from earth a high voltage terminal where they at then **Stripped** a some **electrons** and the resultant positive ions are further accelerated down to earth potential as

- Relative ID across production types and disciplines (by Kullback-Leibler divergence)
  - synchronic: abstracts > research articles; diachronic: t2 > t1

#### Abstract

A growing amount of work has been invested in networks under worst-case scenarios rather than under makes use of the model of "adversarial queuing the (1) (2001) 13–38], under which an adversary is allow

M. Adler, A. Rosén / Journal of Algorithms 55 (2005) 10.

1. Introduction

102

The behavior of packet-switching networks, in which pack network in a continuous manner, has been the subject of conside recent years. See, e.g., [2,3,5–12,14–18]. In such networks packet: adjacent switches over links, in discrete time steps, a prescribed each link in any time step. New packets are injected into the net

### Corpus abs

P(unit | Abs)

*P*(*unit* | RA)

Corpus RA

$$H(Abs) = -\sum_{i} P(unit_{i} | Abs) \log_{2} P(unit_{i} | Abs)$$
  

$$\rightarrow \text{ compare}$$
  

$$H(RA) = -\sum_{i} P(unit_{i} | RA) \log_{2} P(unit_{i} | RA)$$

- Relative ID across production types and disciplines (by Kullback-Leibler divergence) to identify linguistic patterns that might have changed/ are different
  - Use LM trained on one type on a different type
  - Evaluate largest change in log likelihood





Synchronic analyses on Abstracts vs. Research articles

- 1. Classification with SVM
- 2. Entropy rates on single texts
- 3. LMs

# Analysis 1 – SVM Classification Abstracts vs. RAs

Research design

- By possible features involved in reduction and densification
- With Weka (SMO), 10-folds cross-validation, normalized data

type	feature (example)	type	feature (example)
Densification	Sttr Lex. word (nn, adj, adv, vv) / sent.	Modality	Modal verbs ( <i>can, would</i> ) modal meanings (obligation)
<b>Complex NPs</b>	term patterns (adj-n, n-of-n)	Theme	Experiential (NPs)
Simple clauses	X-be-Y		Interpersonal ( <i>Interestingly</i> )
Reduction	Personal pronouns ( <i>we, our</i> ) Definite/indefinite article ( <i>the, a</i> ) Relativizer ( <i>which_that</i> )		Textual ( <i>But, Therefore</i> ) Conj. types at sentence beg. (adversative, additive)
	Affixes ( <i>non-, -like</i> ) Hyphen-words ( <i>degree-3-vertex</i> )	Expansion	Conjunctions ( <i>as, since</i> ), Prepositions ( <i>at, by</i> )

## Analysis 1 – SVM Classification Abstracts vs. RAs

### 70/80s (SASCITEX)

2000s (DASCITEX)



# Analysis 1 – SVM Classification Abstracts vs. RAs

### 70/80s (SASCITEX)

### 2000s (DASCITEX)

### **Typical for abstracts**

type	feature	svm-weight
theme	experiential-sb	-5.03
doncification	sttr	-4.66
uensincation	lex/s	-3.82
datarminara	dt-indef-sb	-4.02
determiners	dt-def-sb	-2.92

type	feature	svm-weight
dancification	sttr	-5.66
densincation	lex/s	-0.88
simple clause	x-be-y	-2.15
determiners	dt-def-sb	-1.06
reduction	pers-pronoun	-0.86

- Densification quite typical
- Simple clause usage
- Reduction to we
   (~60% > than in RAs)



#### Typical for RAs

type	feature	svm-weight
expansion	conjunctions	2.73
simple clause	x-be-y	1.43
determiners	dt-def	1.35
	obligation	1.28
modality	volition	1.08

type	feature	svm-weight
ovnancion	conjuctions	0.97
expansion	noun-relativizer	0.53
modality (	modals	0.90
modality	obligation	0.57
theme	ex-there	0.41

## Analysis 2 – Cross-entropy Abstracts vs RAs

Research design

Based on Genzel and Charniak (2002)

- Procedure designed by Peter Fankhauser
  - Cross-entropy calculated on full articles from 2000s (DASCITEX)
  - Cross-entropy rates observed separately for Abstracts and RAs
    - Headlines not considered (low cross-entropy rates)



## Analysis 2 – Cross-entropy Abstracts vs RAs

• Focus on 5 disciplines from 2000s (DASCITEX)







Bioinformatics





## Analysis 2 – Word Cross-entropy Rates Abstracts-CompSci

ABS-sent1	cross_sent									
Younas2006	15.221	Composite	web	services	provide	promising	prospects	for	conducting	cross-organizational
Tan2006a	11.805	ТСР	over	bandwidth	asymmetric	networks	such	as	Cable	TV
Liao2006	10.65	Authentication	ensures	that	system's	resources	are	not	obtained	fraudulently
Genest2006	10.257	Message	sequence	charts	(	MSC	)	and	High-level	MSC
Pemantle2005a	10.23	Gosper's	algorithm	is	а	cornerstone	of	automated	summation	of
Buhler2005	10.188	Large-scale	comparison	of	genomic	DNA	is	of	fundamental	importance
Cole2006	10.099	We	study	economic	incentives	for	influencing	selfish	behavior	in
Badger2004	10.089	Many	fundamental	questions	in	evolution	remain	unresolved	despite	the
Koch2005	9.902	Graph-based	specification	formalisms	for	access	control	(	AC	)
Lee2005a	9.757	Asynchronous	cellular	automata	(	ACA	)	are	cellular	automata
ABS-sent2	cross_sent			_						
Younas2006	10.045	Such	transactions	:	are	generally	complex	,	require	longer
Tan2006a	5.177	A	number	of	techniques	have	been	proposed	to	address
Liao2006	8.498	Password	authentication	is	one	of	the	simplest	and	the
Genest2006	9.896	They	usually	represent	incomplete	specifications	of	required	or	forbidden
Pemantle2005a	7.461	Milenkovic	and	Compton	in	2002	gave	an	analysis	of
Buhler2005	9.923	То	perform	large	comparisons	efficiently	,	BLAST	(	Methods
Cole2006	6.451	We	consider	а	model	of	selfish	routing	in	which
Badger2004	8.273	This	state	of	affairs	is	partly	due	to	the
Koch2005	8.623	А	security	policy	framework	specifies	а	set	of	(
Lee2005a	8.671	Because	of	the	unpredictability	of	the	order	of	update

## Analysis 2 – Word Cross-entropy Rates Abstracts-CompSci

Further inspection

• By other visualization options

BS-sent1		cross 20-18	cross 17-16	cross 1	L5-10	cross 9-7	cross	6-0
ounas 2006	15.221	3.33	3.64		2.00	0.00		1.25
an2006a	11.805	0.77	2.80		4.00	0.42		2.50
iao2006	10.65	1.54	1.67		3.33	1.82		2.73
ienest2006	10.257	1.85	0.77		1.92	0.40		5.20
emantle200	10.23	0.00	1.82		3.64	1.00		4.00
uhler2005	10.188	1.88	0.67		2.67	0.00		5.00
ole2006	10.099	0.91	2.00		2.00	1.11		4.44
adger2004	10.089 <mark></mark>	0.53	0.00		5.56	1.18		2.94
och2005	9.902	0.61	0.00		4.69	0.65		4.19
ee2005a	9.757	0.00	0.50		4.74	1.11		4.71
BS-sent2		cross 20-18	cross 17-16	cross 1	L5-10	cross 9-7	cross	6-0
. <b>BS-sent2</b> ounas2006	10.045	cross 20-18 0.56	cross 17-16 0.00	cross 1	L5-10 6.25	cross 9-7 1.33	cross	6-0 3.57
<b>.BS-sent2</b> ounas2006 an2006a	10.045 5.177	cross 20-18 0.56 0.00	cross 17-16 0.00 0.00	cross 1	6.25 2.00	cross 9-7 1.33 1.11	cross	6-0 3.57 10.00
<b>BS-sent2</b> ounas2006 an2006a iao2006	10.045 5.177 8.498	cross 20-18 0.56 0.00 0.00	cross 17-16 0.00 0.00 1.25	cross 1	6.25 2.00 2.67	cross 9-7 1.33 1.11 2.86	cross	6-0 3.57 10.00 4.62
<b>BS-sent2</b> ounas2006 an2006a iao2006 ienest2006	10.045 5.177 8.498 9.896	cross 20-18 0.56 0.00 0.00 0.00	cross 17-16 0.00 0.00 1.25 1.54	cross 1	6.25 2.00 2.67 5.83	cross 9-7 1.33 1.11 2.86 0.91	cross	6-0 3.57 10.00 4.62 3.00
<b>BS-sent2</b> ounas2006 an2006a iao2006 ienest2006 emantle200	10.045 5.177 8.498 9.896 7.461	cross 20-18 0.56 0.00 0.00 0.00 0.00	cross 17-16 0.00 0.00 1.25 1.54 0.00		6.25 2.00 2.67 5.83 4.74	cross 9-7 1.33 1.11 2.86 0.91 2.22	cross	6-0 3.57 10.00 4.62 3.00 4.12
<b>BS-sent2</b> ounas2006 an2006a iao2006 ienest2006 emantle200 uhler2005	10.045 5.177 8.498 9.896 7.461 9.923	cross 20-18 0.56 0.00 0.00 0.00 0.00 1.36	cross 17-16 0.00 1.25 1.54 0.00 0.48		6.25 2.00 2.67 5.83 4.74 2.50	cross 9-7 1.33 1.11 2.86 0.91 2.22 2.63		6-0 3.57 10.00 4.62 3.00 4.12 3.89
BS-sent2 ounas2006 an2006a iao2006 Genest2006 emantle200 uhler2005 ole2006	10.045 5.177 8.498 9.896 7.461 9.923 6.451	cross 20-18 0.56 0.00 0.00 0.00 1.36 0.24	cross 17-16 0.00 1.25 1.54 0.00 0.48 0.24		6.25 2.00 2.67 5.83 4.74 2.50 1.50	cross 9-7 1.33 1.11 2.86 0.91 2.22 2.63 1.54	Cross	6-0 3.57 10.00 4.62 3.00 4.12 3.89 7.11
<b>BS-sent2</b> ounas2006 an2006a iao2006 ienest2006 emantle200 uhler2005 ole2006 adger2004	10.045 5.177 8.498 9.896 7.461 9.923 6.451 8.273	cross 20-18 0.56 0.00 0.00 0.00 1.36 0.24 0.45	cross 17-16 0.00 1.25 1.54 0.00 0.48 0.24 0.00		6.25 2.00 2.67 5.83 4.74 2.50 1.50 3.50	cross 9-7 1.33 1.11 2.86 0.91 2.22 2.63 1.54 2.11		6-0 3.57 10.00 4.62 3.00 4.12 3.89 7.11 5.00
<b>BS-sent2</b> ounas2006 an2006a iao2006 ienest2006 emantle200 uhler2005 ole2006 adger2004 och2005	10.045 5.177 8.498 9.896 7.461 9.923 6.451 8.273 8.623	cross 20-18 0.56 0.00 0.00 0.00 1.36 0.24 0.45 0.59	cross 17-16 0.00 1.25 1.54 0.00 0.48 0.24 0.00 0.30		5-10 6.25 2.00 2.67 5.83 4.74 2.50 1.50 3.50 2.19	cross 9-7 1.33 1.11 2.86 0.91 2.22 2.63 1.54 2.11 1.94		6-0 3.57 10.00 4.62 3.00 4.12 3.89 7.11 5.00 5.67
<b>BS-sent2</b> ounas2006 an2006a iao2006 ienest2006 emantle200 uhler2005 ole2006 adger2004 och2005 ee2005a	10.045 5.177 8.498 9.896 7.461 9.923 6.451 8.273 8.623 8.671	cross 20-18 0.56 0.00 0.00 0.00 1.36 0.24 0.45 0.59 0.79	cross 17-16 0.00 1.25 1.54 0.00 0.48 0.24 0.00 0.30 0.27		5-10 6.25 2.00 2.67 5.83 4.74 2.50 1.50 3.50 2.19 3.89	cross 9-7 1.33 1.11 2.86 0.91 2.22 2.63 1.54 2.11 1.94 0.86		6-0 3.57 10.00 4.62 3.00 4.12 3.89 7.11 5.00 5.67 4.71

## Analysis 3 LMs on Abstracts vs RAs

Andrews2004

<sup>1</sup>/<sub>4</sub> of each discipline for equal distribution

(sent 18) ...

Research design
LMs build based on pos trigrams (done by Jonathan Poitz based on his BSc-thesis)
LMs build based on word trigrams (done by Anna Currey)

**DASCITEX Corpus** Abstracts RAs B2 B3 B4 C1 C2 C3 C1 C2 C3 C4 B4 B2 B3 B1 Α B1 Α Aggarwal2006 Aggarwal2006 (sent 10) (sent 10)

Andrews2004

<sup>1</sup>/<sub>4</sub> of each discipline for equal distribution

(sent 18) ...

Training cot

Set / Instances

Training set ~ 2000 each

C4

Test set ~ 500 each

## Analysis 3 LMs on Abstracts vs RAs

Use of LMs

- 1. Train model on Abstracts/RAs and
  - a. Test how well the models distinguish between the two (similar to a classification task)
     → are they distinct in terms of Inf Theory
  - b. Compare relative ID of Abstracts and RAs to determine which have a higher ID
     → which is more informationally dense
- Train model on Abstracts and test on RAs and train model on RAs and test on Abstracts
   → to identify linguistic patterns that are different

## Analysis 3 LMs on Abstracts vs RAs

Preliminary analysis (done by Jonathan Poitz)

- Tested LM\_Abs and LM\_RAs on Abstracts/RAs
- Lower perplexity shows the right "class assignment" and a strength of the strengt

abstracts	PPI	ABS PP	L RAs	diff	RAs	PPL ABS PP	L RAs	diff
Aaltonen2007_C3	₽	7.26 🔶	8.29	1.03	Aaltonen2007_C3	1.26 🔶	8.26	0.99
Abney2004_B1	₽	9.53 🔶	10.35	0.82	Abney2004_B1	🔶 12.72 🖊	10.41	2.31
Adachi2006_C4	₽	5.93 🔶	6.69	0.77	Adachi2006_C4	1.02 🕂	7.89	1.13
Adger2005_C1	₽	7.82 🔶	8.11	0.29	Adger2005_C1	🔶 11.65 🖊	9.89	1.76
Aggarwal2006_A	₽	7.12 🔶	7.89	0.76	Aggarwal2006_A	🔶 10.17 🖊	8.89	1.28
Aguilar2005_B3	$\mathbf{I}$	6.29 🔶	8.32	2.03	Aguilar2005_B3	🔶 9.98 🖊	8.73	1.25

	class	Abstracts	RAs
	abs_A	46	10
	abs_B1	51	5
	abs_B2	56	0
,,	abs_B3	53	3
	abs_B4	56	0
	abs_C1	43	12
	abs_C2	54	1
	abs_C3	53	2
	abs_C4	55	0
	RAs_A	0	56
	RAs_B1	0	56
	RAs_B2	0	56
	RAs_B3	0	56
	RAs_B4	0	56
	RAs_C1	0	55
	RAs_C2	0	55
	RAs_C3	0	55
	RAs_C4	0	55

- → Why 100% for RAs?
- First analyses are always very valuable!
- Help to reconsider research design!

For example:

- Equal size of Abstracts and RAs
- Equal proportion of disciplines
- Similar data composition (include/exclude headlines, etc.)

## Conclusion

- First analyses showed that Abstracts differ in terms of information density from RAs
- BUT: Research design and interpretation of the data is not trivial!
  - Have a good knowledge of our data
  - Gain a better understanding of what the methods can provide us with to ask the right question that can be answered by the appropriate methodology
  - Need of good visualizations
  - Etc.

## Thank you for your attention!

