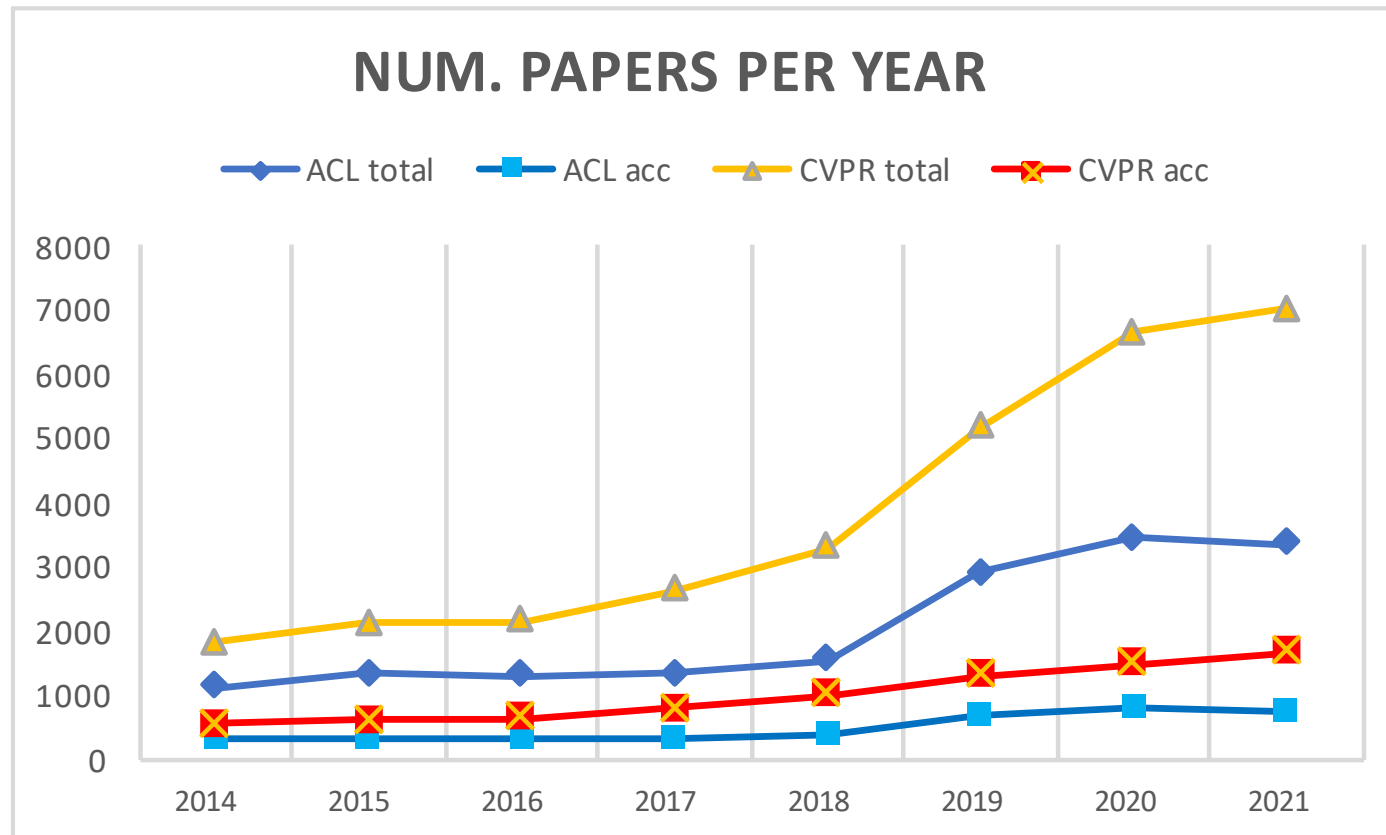


What do writing features tell us about AI papers?

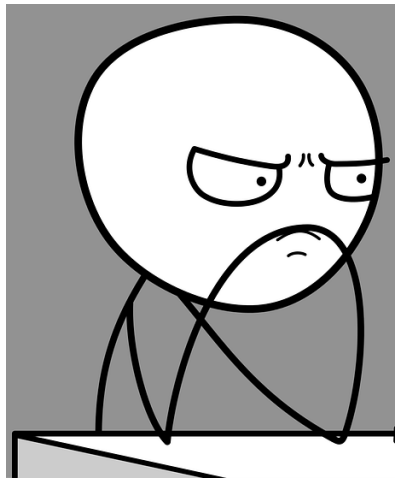
Zining Zhu, Bai Li, Yang Xu, Frank Rudzicz

Recent submissions increase in numbers



Data source: <https://github.com/lixin4ever/Conference-Acceptance-Rate>

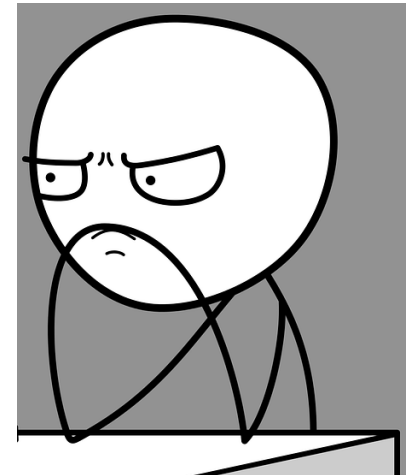
Problems from two sides



- Poorly organized
- Methodology is problematic
- Result is unclear
- Question - analysis mismatch
- Limited novelty
- Limited impact
- Ethical concerns

Gets random submissions

- Didn't read carefully
- Doesn't understand our method
- Doesn't think hard
- Doesn't understand the field
- Reviewers are paranoid



Gets random peer reviews

Possibilities of improvements?

- Improved peer review procedure
 - OpenReview
 - ACL Rolling Review
- Use DNN to predict paper outcomes
 - Text classification problem
- Intuition: **text markers** can lead to scalable solutions
 - “Best of both worlds”
 - Similar: Automatic Essay Scoring, e.g., Grammarly →

Overall score **96**
See performance

Goals
3 of 5 set


All suggestions

Correctness
1 alert

Clarity 
Very clear

Engagement
Engaging

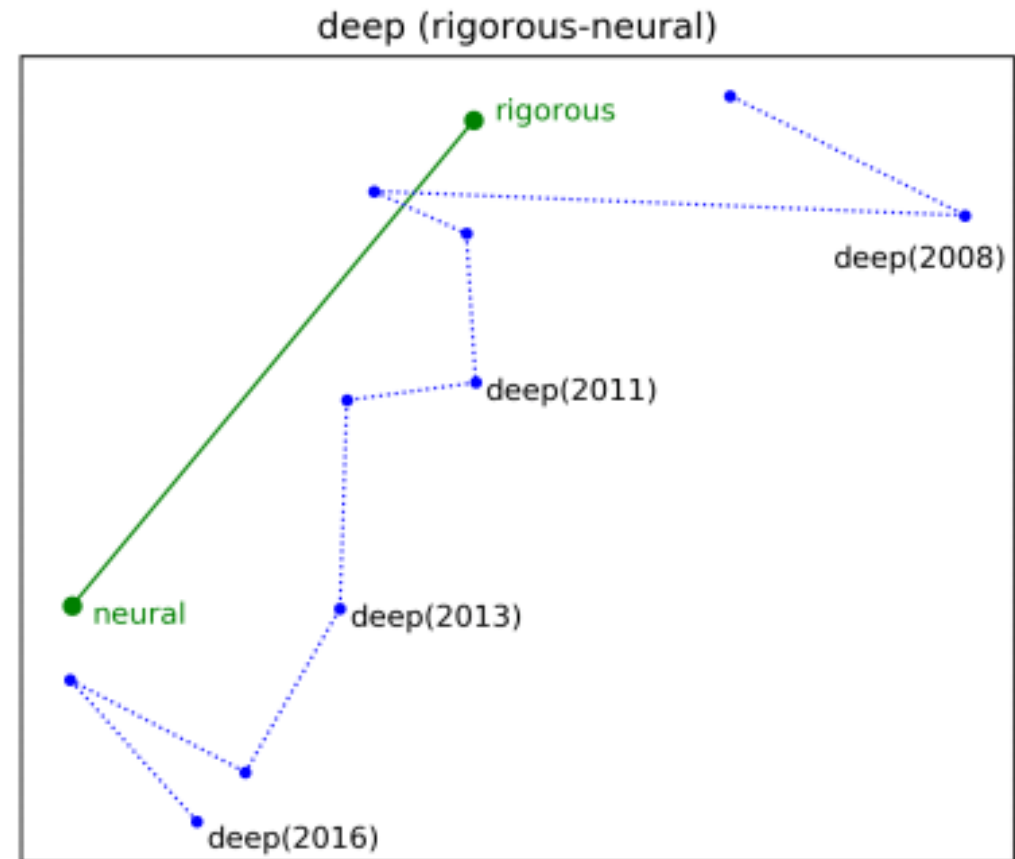
Delivery 
Just right

Premium 
Advanced suggestions

There are some interesting text markers

An example: locations on the *semantic coordinates*.

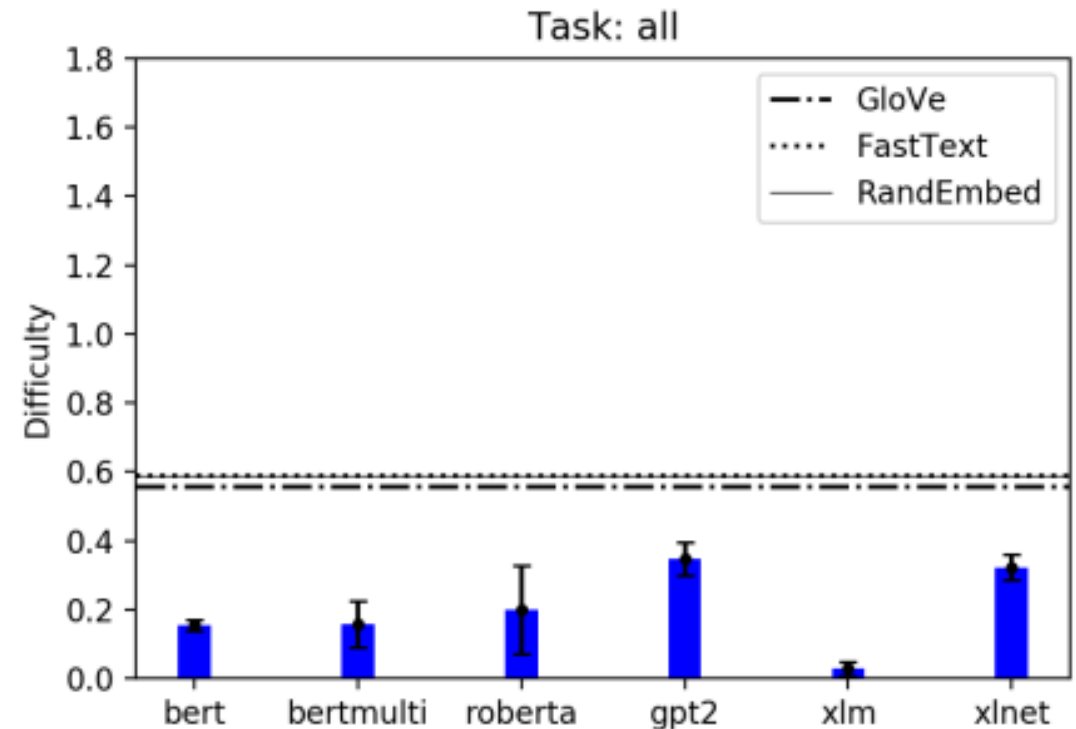
- Hypothesis: word semantics shift along certain coordinates.
- Semantically stable words form **coordinates**.
- **Target words** shift along the coordinates.



And they can be useful

An example: examine the rhetorical capacities of neural LMs.

- Use simple models (“probes”) to predict rhetorical features
- Use loss to measure the difficulty → negation of “the goodness of encoding the knowledge”



Let's try some text markers for AI papers

We consider *74 writing features*

i.e., do **not** explicitly describe the semantics.

- Metadata: outbound citations, article lengths, sentence lengths...
- Readability: Flesch, Flesch-Kincaid, semantic surprisal
- Lexical richness: Moving-average type-token ratio
- Syntactic: Grammar error counts, active / passive voice portions
- Stylistic features: POS signal constituency, RST signal constituency

What can the writing features do?

- They are **correlated to** Conference (C) vs Workshop (W) appearance.

Venue	Features	Spearman R	ATE	Interpretation
ACL	flesch_kincaid_grade_level_bodytext	-0.05	+0.05	Ambiguous
	<u>grammar_errors_abstract</u>	-0.09**	-0.01	W papers are larger
	surprisal_abstract_std	-0.01	+0.00	Ambiguous
	title_word_length	-0.09**	-0.01	W papers are larger
	voice_bodytext_active	+0.09**	+0.15	C papers are larger
EMNLP	outbound_citations_per_word	-0.17**	+67.6	Ambiguous
	n_author	-0.17**	-0.05	W papers are larger
	<u>grammar_errors_abstract</u>	-0.18**	+0.01	W papers are larger
	n_outbound_citations	-0.09	+0.09	Ambiguous
	abstract_word_counts	-0.16**	+0.00	W papers are larger

What can the writing features do?

- They can predict Conference (C) vs Workshop (W) appearance.

Venue Name	74 features	Writing Features					TF-IDF		RoBERTa
		RST	Surprisal	Grammar	LexRich	Readability	Full text	Abstract	Abstract
AAAI	.755(.028)	+0.001	+0.024	+0.010	-0.002	+0.009	+0.206**	+0.203**	+0.212**
ACL	.867(.004)	+0.001	+0.000	+0.001	+0.001	+0.001	-0.004	-0.008*	-0.015
COLING	.837(.010)	+0.005	+0.005	+0.003	+0.003	+0.004	+0.049**	+0.051**	+0.052**
CVPR	.900(.005)	-0.007	-0.006	-0.001	-0.006	-0.005	-0.052**	-0.067**	-0.070**
EMNLP	.737(.020)	+0.003	+0.014	+0.012	+0.022	+0.015	+0.159**	+0.153**	+0.102
ICML	.659(.023)	-0.277**	-0.042*	-0.102	-0.262**	-0.185	+0.333**	+0.333**	+0.300**
IJCAI	.868(.002)	-0.067**	-0.066**	-0.045**	+0.075**	-0.067**	-0.029**	-0.061**	-0.076
NAACL	.757(.019)	+0.016	+0.016	+0.011	+0.017*	+0.016*	-0.107**	-0.128**	-0.182
NeurIPS	.586(.035)	-0.193**	-0.039	-0.097	-0.212**	-0.157	+0.031	-0.077*	-0.110

Table 3: F1 scores of the C vs. W classification results. The second column shows the F1 score using 74 writing features. The remaining columns show the values relative to the baseline. * $p < .005$ and ** $p < .001$ respectively, both on 2-tailed t -test with $dof = 10$, Bonferroni correction.

Sometimes comparable to TF-IDF features, and even RoBERTa.

What can the writing features do?

- They can **sort of** tell apart between different venues.

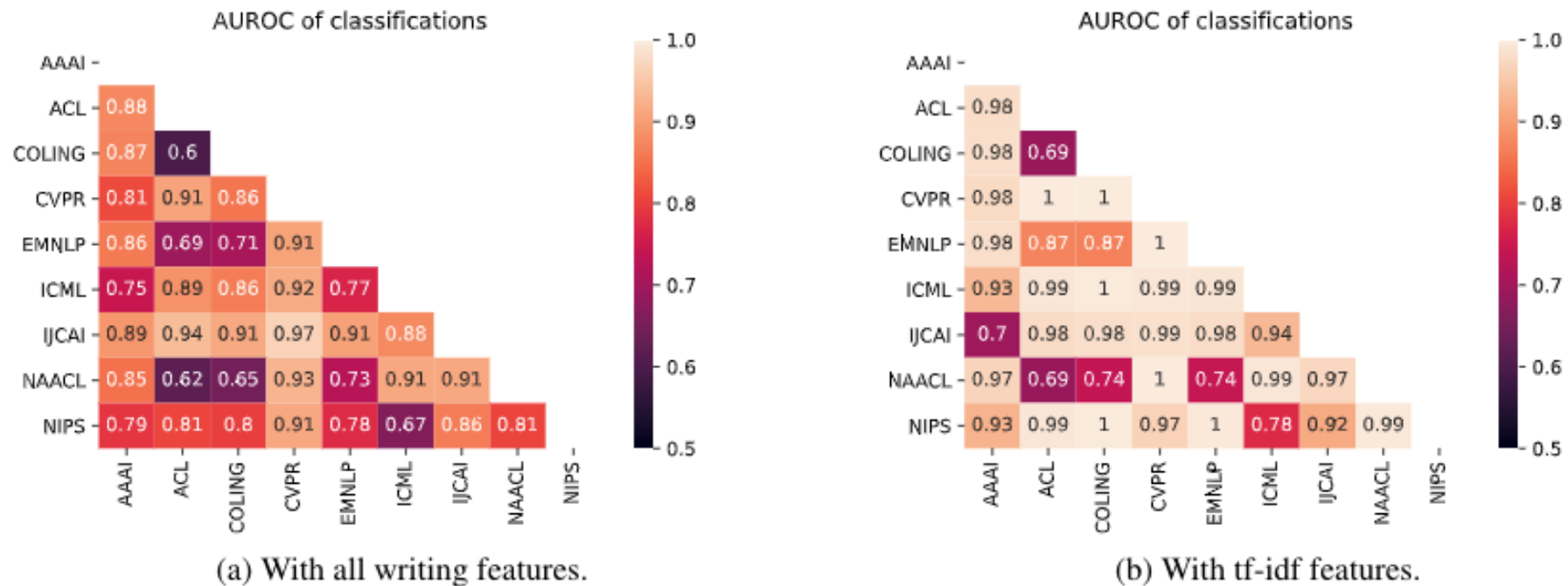


Figure 1: The AUROC of inter-venue classifications. The venues in the same categories (e.g., COLING and ACL) are harder to tell apart than other venues, using either the writing features or tf-idf features.

What can the writing features do?

- They **can not** predict the inbound citation counts.

Venue Name	Writing Features						Baseline	TF-IDF	
	74 features	RST	Surprisal	Grammar	LexRich	Readability		Full text	Abstract
AAAI	+8.61	+0.09	+0.69	+0.27	+1.16	+0.61	20.77(27)	0.07(.02)	0.07(.01)
ACL	+6.27	+0.48	+0.11	+0.10	+0.34	+2.89	389.76(636)	0.15(.01)	0.17(.01)
COLING	+248.87	+3.51	+0.05	+0.32	+0.62	+368.18	437.76(1006)	0.15(.02)	0.16(.01)
CVPR	-6.11	+239.82	+9.62	+6.90	+22.24	+12.35	15273.45(24710)	0.17(.01)	0.19(.01)
EMNLP	+55211.48	+8.12	+4.66	+42.06	+12.11	+458.65	1194.59(2788)	0.15(.02)	0.17(.03)
ICML	+45.13	+6.93	+37.59	+9.77	+988.27	+86.97	1279.15(1200)	0.02(.01)	0.02(.02)
IJCAI	+8.84	+1.37	+0.81	+1.20	+3.77	+1.97	23.77(25)	0.16(.01)	0.22(.04)
NAACL	+18.04	+2.22	+0.83	+0.99	+0.08	+195.92	420.34(855)	0.22(.01)	0.22(.01)
NeurIPS	+78.25	+4.64	+6.37	+39.81	-2.87	+76.99	3305.99(5216)	0.20(.01)	0.23(.01)

But TF-IDF features
can predict!

More about the data

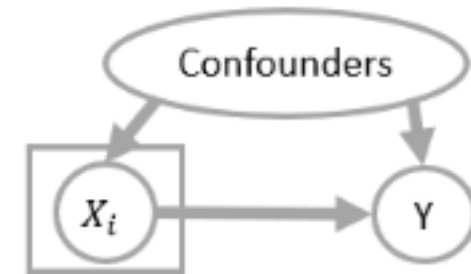
- Computed features on 945,674 CompSci articles from S2ORC.
 - 97.68% have ≤ 10 annual income citations.
 - Each article is cited 1.59 (std=13.5) times per year.
- Gave C & W labels for AI venues.
 - NLP: ACL, COLING, EMNLP, NAACL
 - AI: AAAI, IJCAI
 - ML: ICML, NeurIPS
 - CV: CVPR
 - ICRA and ICASSP not used

Venue Name	N. articles	N. articles by label	
		C	W
AAAI	624	395	229
ACL	2,836	2,175	661
COLING	1,860	1,353	507
CVPR	3,495	2,824	671
EMNLP	714	437	277
ICML	930	396	534
ICRA	703	662	41
IJCAI	632	423	209
NAACL	2,142	1,354	788
NeurIPS	930	396	534

Table 7: Number of C and W articles of each venue. The arXiv papers of the corresponding sections are included as W papers. For example, `cs.Learning` and `cs.ML` are included in the W portions of ICML and NeurIPS.

More about the writing features...

- They are mutually dependent
 - Causal model assumed independence -> Observe multicollinearity effect.
 - Partial features can often predict well.
- They describe more than “just the writing”.
 - E.g., RST: stylistic choices -> author -> content
 - E.g., title length -> scope of content -> num. readers -> citation counts
- BTW: Good papers are more than well-“written”.
 - Should consider their [impact](#).



Summary

- Computed *74 writing features*
- Compiled a test suite to assess their usefulness:
 - Conference vs. Workshop appearance prediction
 - Venue appearance prediction
 - Citation counts prediction
- Text markers can lead to scalable, high-quality, and trustworthy solutions for assessing academic article writing.
 - More text markers, and group them together.
 - Additional subjects, more than just CompSci / AI

Connections beyond academic writing

Structural probe, edge
probing, and rediscovering
classic NLP pipeline

Diagnostic
classification

Linguistic
theory

How is BERT surprised? (Li et al., 2021)

What does BERT learn about the structure of language? (Jawahar et al 2019)

Interpretable
text-markers

Symbolic
reasoning

ML4H

Predicting the inductive bias (Lovering et al., 2021)

RoBERTa acquires linguistic features eventually (Warstadt et al., 2020)

Lexical features are more vulnerable (Balagopalan et al., 2019)

Agreeing on interpretations of linguistic features (Zhu et al., 2019)