

# *JAMDEC: Unsupervised Authorship Obfuscation using Constrained Decoding over Small Language Models*

---

Jillian Fisher<sup>1</sup>, Ximing Lu<sup>2,3</sup>, Jaehun Jung<sup>2</sup>, Liwei Jiang<sup>2,3</sup>, Zaid Harchaoui<sup>1</sup>, Yejin Choi<sup>2,3</sup>

<sup>1</sup>Department of Statistics, University of Washington, <sup>2</sup>Paul G. Allen School of Computer Science & Engineering, University of Washington, <sup>3</sup>Allen Institute for Artificial Intelligence

# Authorship Obfuscation

Task of altering a text to reduce the discovery of the author's identity.

***Many settings require obscuring authorship...***

Blind Review for  
Scientific Papers

Interaction on Mental  
Health Forums

Anonymous Online  
Review

***Current authorship obfuscation methods have challenges....***

Require Large  
Computing

Involve Proprietary  
LLMs

Require an Extra  
Authorship Corpus



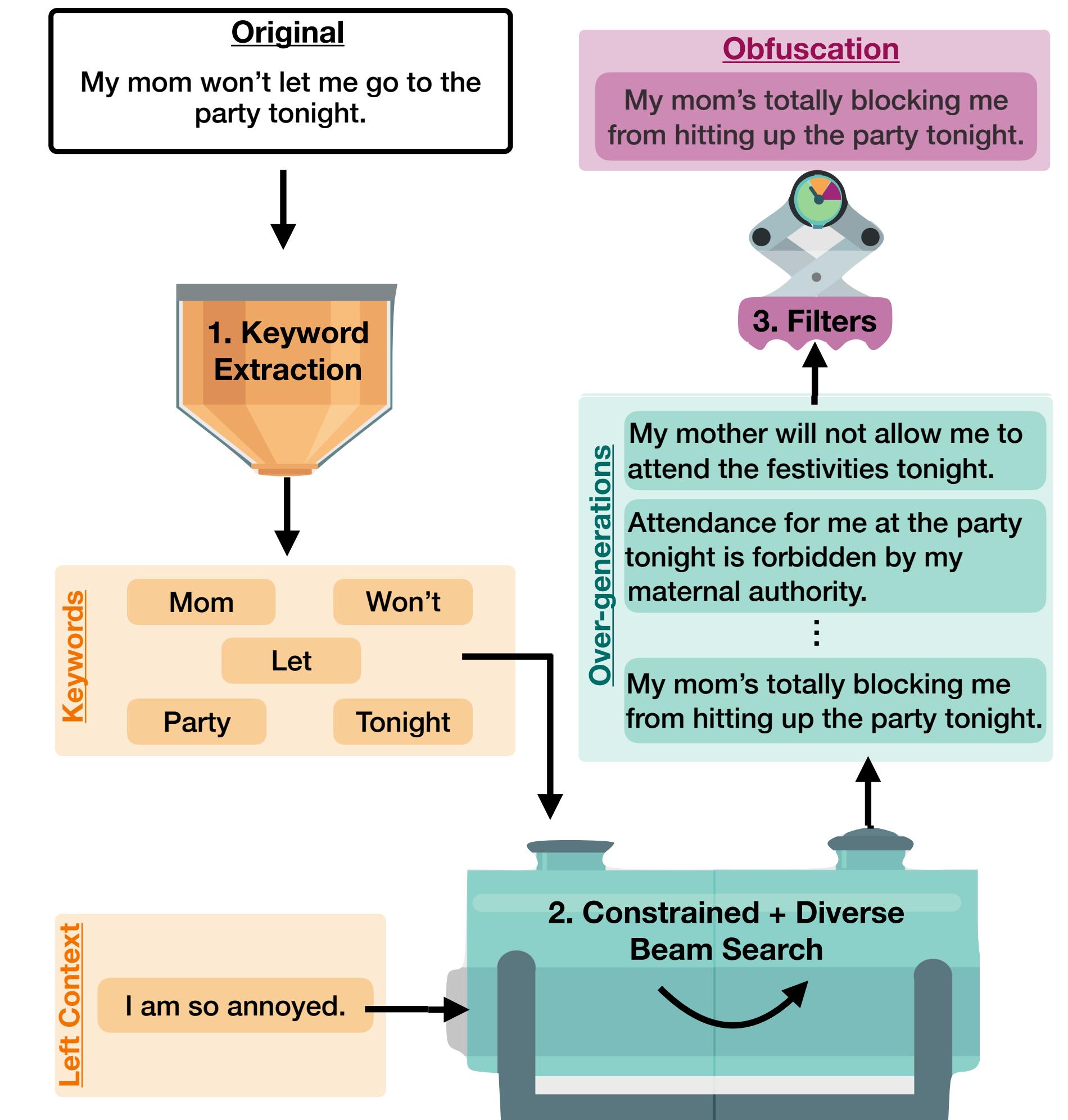
# JAMDEC Decoding

## ***JAMDEC Decoding***

- user-controlled, inference-time algorithm for authorship obfuscation that can be applied to any text and authorship without a separate authorship corpus

## **• 3 Stage Approach:**

1. *Keyword Extraction*: Extract keywords to maintain original content
2. *Constrained + Diverse Beam Search*: Augmented decoding strategy which encourages diverse but constrained generations
3. *Filters*: Maintain fluency and content preservation, +any user-specified control



# JAMDEC Decoding: Innovation

## Keyword Extraction: Likelihood-based Method

- we select the top-k tokens with the lowest conditional probabilities, as measured by a specific language model, as keywords for a given sentence

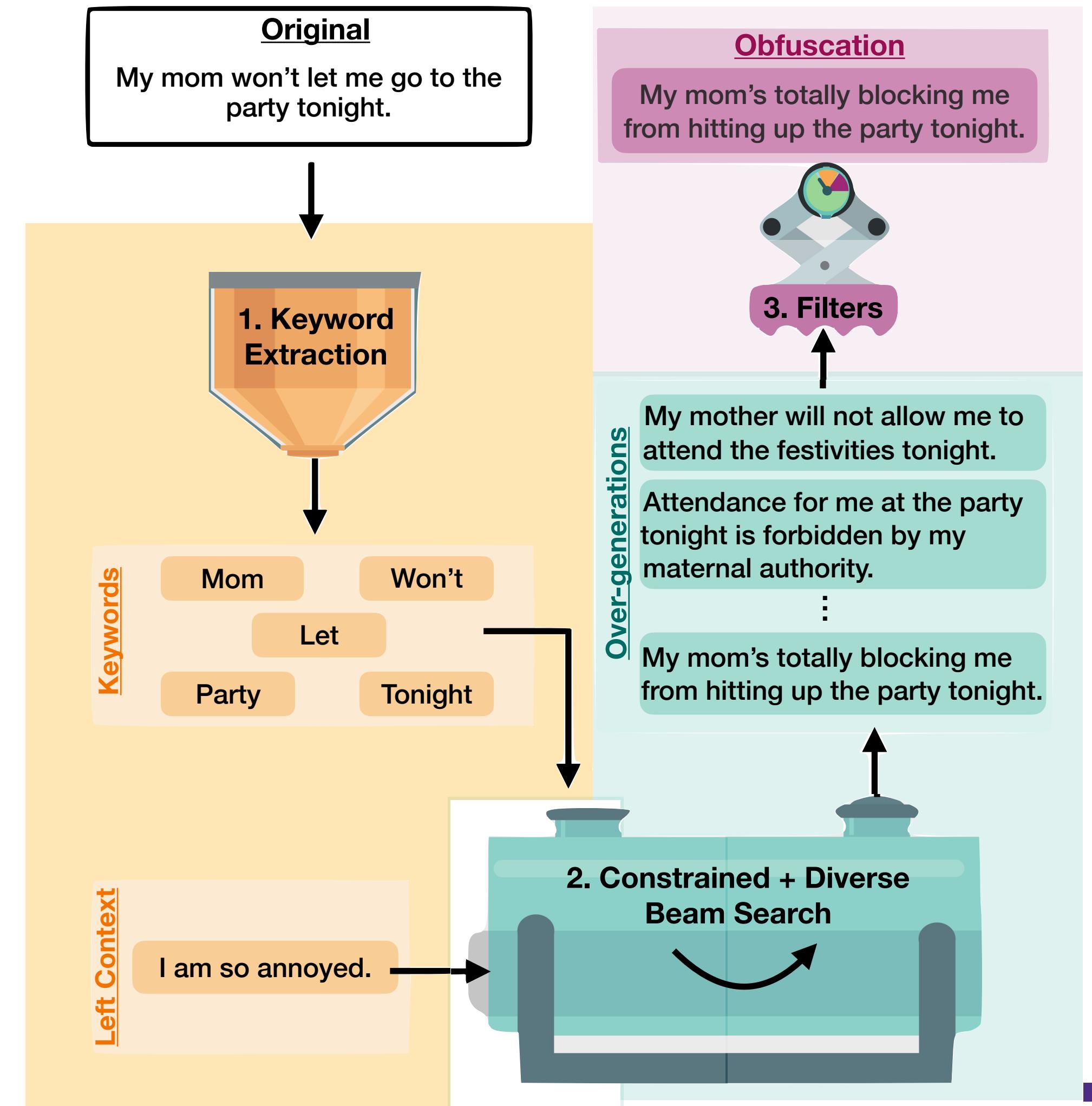
## Constrained + Diverse Beam Search (over-generation)

- Constrained Beam Search is base algorithm, but uses the scoring function from Diverse Beam Search instead of likelihoods when iteratively selecting the top-k candidates from each bank.

$$\arg \max_{w \in W} P_w(y | x) + \lambda_1 D(y, Y) + \lambda_2 C(y)$$

## Filtering

- Reduce pool and allow personalization of user



# JAMDEC Decoding: Results

## ***JAMDEC Decoding Experiments***

- Two Datasets
  - AMT: formal scholarly passages
  - BLOG: diary-style entries
  - Number of Authors: 3, 5, or 10
- Baselines
  - Mutant-X: Iteratively re-writes and combines randomly
  - Paraphrase
  - Round Trip Machine Translation
  - Stylometric: simple changes such as synonyms, number of words, punctuation, etc.
- Metrics
  - Obfuscation (Drop Rate: drop in accuracy compared to original)
  - Fluency: METEOR and CoLA
  - Content Preservation: NLI
  - Overall Metric: Task Score (average of Drop Rate, NLI, and CoLA)

Dataset	Method	Mutant-X		Paraphrase	Machine Transl.	Stylometric	JAMDEC	
		ENS	RFC				W/O Stylo	W/Stylo
AMT-3	Drop Rate (ENS)	*	-0.04	<b>0.04</b>	<b>0.04</b>	-0.03	<b>0.11</b>	<b>0.11</b>
	Drop Rate (BertAA)	<u>0.10</u>	0.04	0.04	0.08	<b>0.12</b>	0.04	0.04
	METEOR	<u>0.80</u>	<b>0.81</b>	0.55	0.69	<b>0.80</b>	0.62	0.62
	NLI	0.60	0.61	0.62	<b>0.75</b>	0.50	<b>0.75</b>	<b>0.81</b>
	CoLA	0.50	0.51	0.78	0.69	0.46	<b>0.85</b>	<u>0.79</u>
	Task Score (ENS)	*	0.36	0.48	<b>0.49</b>	0.31	<b>0.57</b>	<b>0.57</b>
	Task Score (BertAA)	0.40	0.39	0.48	<u>0.51</u>	0.36	<b>0.55</b>	<b>0.55</b>
AMT-5	Drop Rate (ENS)	*	0.08	<b>0.20</b>	<b>0.20</b>	<b>0.23</b>	0.10	0.13
	Drop Rate (BertAA)	<u>0.07</u>	0.00	-0.06	<u>0.07</u>	0.04	<b>0.14</b>	<b>0.14</b>
	METEOR	<u>0.74</u>	0.72	0.57	0.68	<b>0.79</b>	0.61	0.61
	NLI	0.56	0.57	0.62	0.74	0.48	<u>0.76</u>	<b>0.82</b>
	CoLA	0.51	0.55	0.77	0.69	0.46	<b>0.85</b>	<u>0.79</u>
	Task Score (ENS)	*	0.40	0.53	0.54	0.39	<b>0.57</b>	<u>0.58</u>
	Task Score (BertAA)	0.38	0.37	0.44	<u>0.50</u>	0.33	<b>0.58</b>	<b>0.58</b>
AMT-10	Drop Rate (ENS)	*	0.10	0.07	0.19	0.11	<b>0.44</b>	<u>0.41</u>
	Drop Rate (BertAA)	0.03	<u>0.04</u>	-0.04	<b>0.06</b>	0.00	-0.03	-0.02
	METEOR	<u>0.84</u>	<b>0.86</b>	0.54	0.66	0.81	0.60	0.61
	NLI	0.61	0.64	0.61	<b>0.73</b>	0.45	<b>0.79</b>	<b>0.79</b>
	CoLA	0.53	0.57	<u>0.77</u>	0.68	0.46	<b>0.78</b>	<b>0.78</b>
	Task Score (ENS)	*	0.44	0.48	0.53	0.34	<b>0.67</b>	<u>0.66</u>
	Task Score (BertAA)	0.39	0.42	0.45	0.49	0.30	<u>0.51</u>	<b>0.52</b>

\* Similar results for the BLOG dataset (see paper for all results)

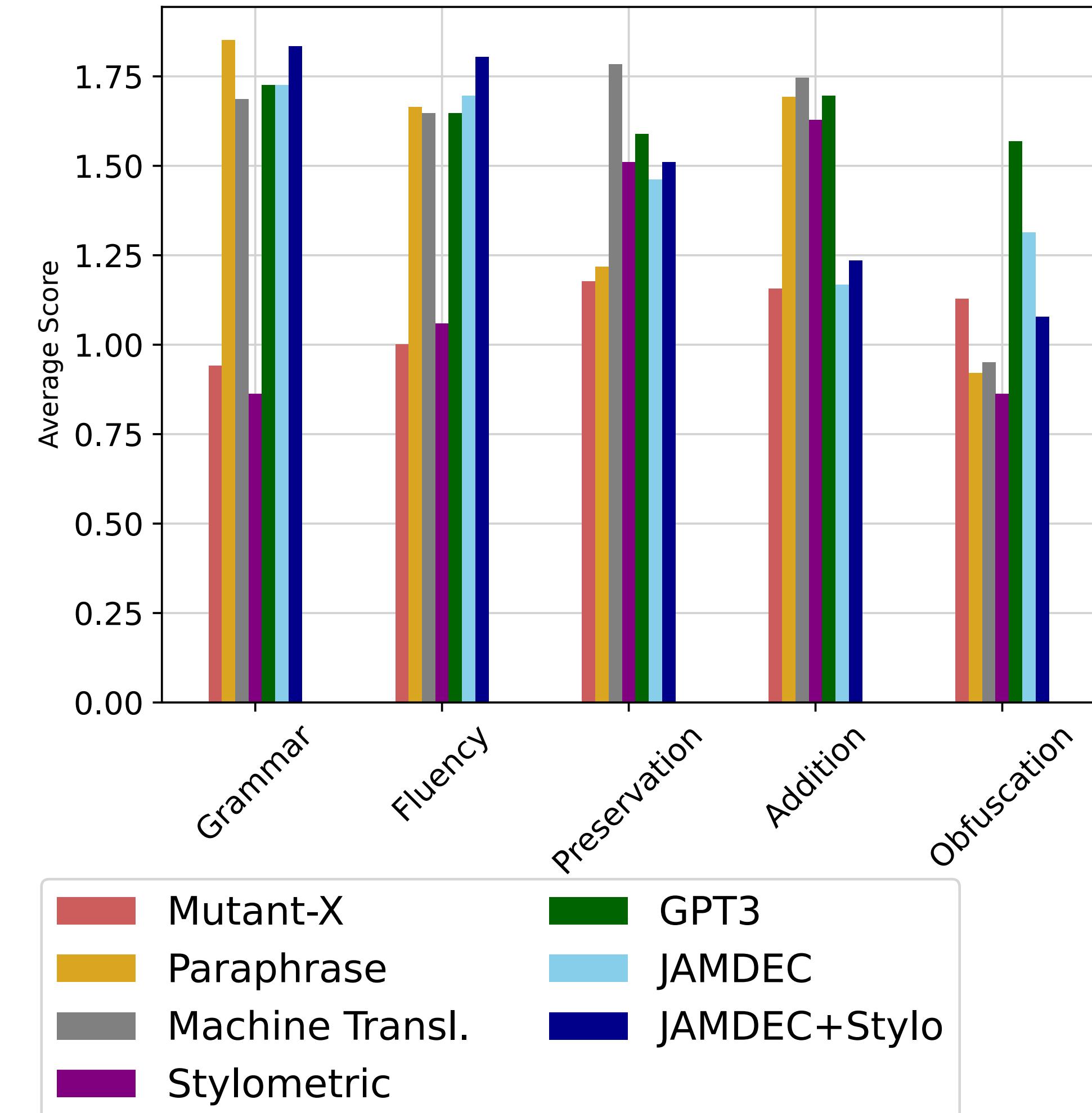


# JAMDEC Decoding: Results

## ***JAMDEC Decoding Experiments***

- Two Datasets
  - AMT: formal scholarly passages
  - BLOG: diary-style entries
  - Number of Authors: 3, 5, or 10
- Baselines
  - Mutant-X: Iteratively re-writes and combines randomly
  - Paraphrase
  - Round Trip Machine Translation
  - Stylometric: simple changes such as synonyms, number of words, punctuation, etc.
- Metrics
  - Human Evaluation

**Human Evaluation**

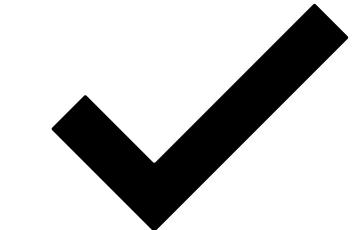
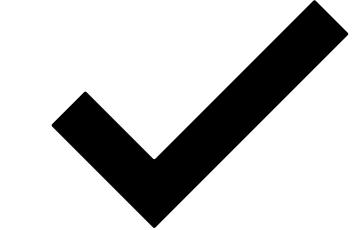


# JAMDEC Decoding: Results

*Performs similar to much larger models!*

Method	GPT3.5		MASQUERADE	
	Sentence	Paragraph	W/O Stylo	W/ Stylo
Metric				
Drop Rate (ENS)	<b>0.23</b>	<b>0.23</b>	<u>0.11</u>	<u>0.11</u>
Drop Rate (BertAA)	<b>0.13</b>	<u>0.09</u>	0.04	0.04
METEOR	0.33	<u>0.41</u>	<b>0.62</b>	<b>0.62</b>
NLI	<u>0.77</u>	0.73	0.75	<b>0.81</b>
CoLA	0.76	<u>0.80</u>	<b>0.85</b>	0.79
Task Score (ENS)	<b>0.59</b>	<b>0.59</b>	<u>0.57</u>	<u>0.57</u>
Task Score (BertAA)	<b>0.55</b>	<u>0.54</u>	<b>0.55</b>	<b>0.55</b>

# JAMDEC Decoding: Results

Method	Generation	
Original	The Ex. An ex holding a grudge can do a lot of damage in a short amount of time. He knows enough to open accounts in your name, and he has the motive to hurt you.	
Mutant-X	The Ex. An ex holding a <b>bitterness able ought</b> a lot of damage in a <b>length quantity</b> of time. He knows enough to <b>ascend</b> accounts in <b>Your prefix</b> , and he has the <b>justifiable to impair You</b> .	Ungrammatical
Paraphrase	<b>A lot of damage can be done In a short period of time.</b> He knows <b>how to</b> open accounts In your name and he <b>wants</b> to hurt you.	Incorrect Content
Machine Translation	<b>The former.</b> An <b>old man who holds a knife</b> can make a lot of damage in a short time. He knows enough to open accounts in your name, and he has the <b>reason</b> to hurt you.	Incorrect Content
Stylometric	An ex <b>holding, a</b> grudge can do a lot <b>inside damage</b> in a <b>brief</b> amount in time, <b>yet</b> he knows enough to open accounts in your name, and he has the motive to hurt you.	Missing Meaning
JAMDEC	The Ex. <b>When the ex is holding his grudge against the person who caused him lot of damage to his life, he is short sighted and will do anything in his power to get back at that person, no matter how much it will hurt the person he is trying to get revenge against.</b> He knows enough to open accounts in your name, and he has the motive to hurt you.	
JAMDEC + Stylo	The Ex. <b>When the ex is holding his grudge against the person who caused him lot of damage to his life, he is short sighted and will do anything in his power to get back at that person, no matter how much it will hurt the person he is trying to get revenge against.</b> He <b>believes</b> enough to open accounts in your name, and he has the <b>reason</b> to hurt you.	

# More in the Paper

- Comparison of trade-off between obfuscation, content-preservation, and grammaticality
- Ablation of JAMDEC Method (different beam width, with/without diversity, different filters, etc.)
- Comparison of “Style Transfer” methods
- Evaluation using “Adversarial Threat Models”
- Time Consumption Analysis (it’s competitive!)
- Discussion of similarity to other tasks (paraphrasing, style transfer, authorship attribution, etc.)
- *And MORE!*

# Thank You!

Paper



Code



W