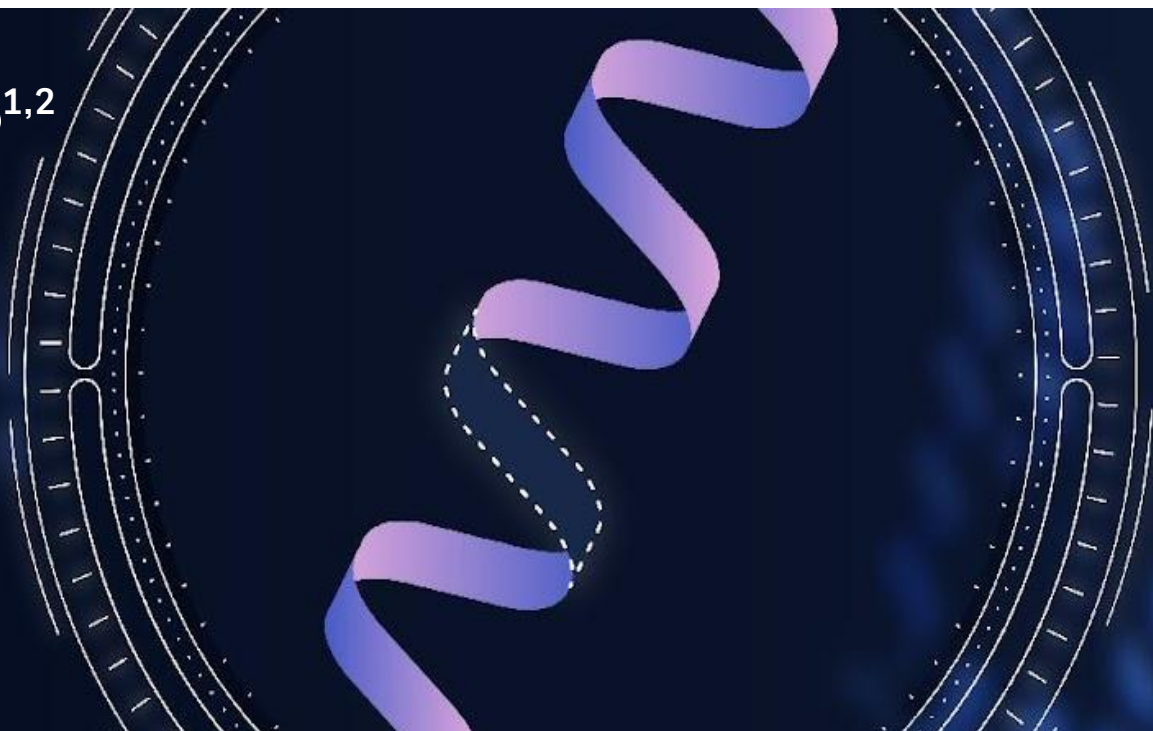


Efficient semantic uncertainty quantification in language models via diversity-steered sampling

Ji Won Park¹, Kyunghyun Cho^{1,2}

¹ Prescient Design, Genentech

² Center for Data Science, NYU



Estimating uncertainty in language generation is critical for building reliable AI systems but often relies on drawing large IID sample sets.

Q: What is a way to measure the Hubble constant?

A: Use **Cepheid variable stars**.

Cluster 1

A: By analyzing the distances to **Cepheid variable stars**.

A: By comparing the brightness and redshift of Type Ia **supernovae**.

Cluster 2

A: By measuring the distances to Type Ia **supernovae** in different galaxies.

A: The Hubble constant can be measured using the light curves of Type Ia **supernovae**, which have a constant maximum brightness allowing for distance estimation.

A: The method of standard candles involves using **supernovae** as "candles" to measure distances and then infer the Hubble constant from redshift data.

A: One way to measure the Hubble constant is by using the distance-luminosity relationship for Type Ia **supernovae**.

For example, **semantic entropy** (Kuhn et al., 2023) groups generations into semantically equivalent clusters and computes the entropy over the clusters.

Q: What is a way to measure the Hubble constant?

A: Use **Cepheid variable stars**.

Cluster 1

A: By analyzing the distances to **Cepheid variable stars**.

$$p(c | x, \theta) = \int 1[y \in c | x, \theta] p(y | x, \theta) dy$$

$$H(p(c | x, \theta)) = - \sum_{c \in \mathcal{C}} p(c | x, \theta) \log p(c | x, \theta)$$

A: By comparing the brightness and redshift of Type Ia **supernovae**.

Cluster 2

A: By measuring the distances to Type Ia **supernovae** in different galaxies.

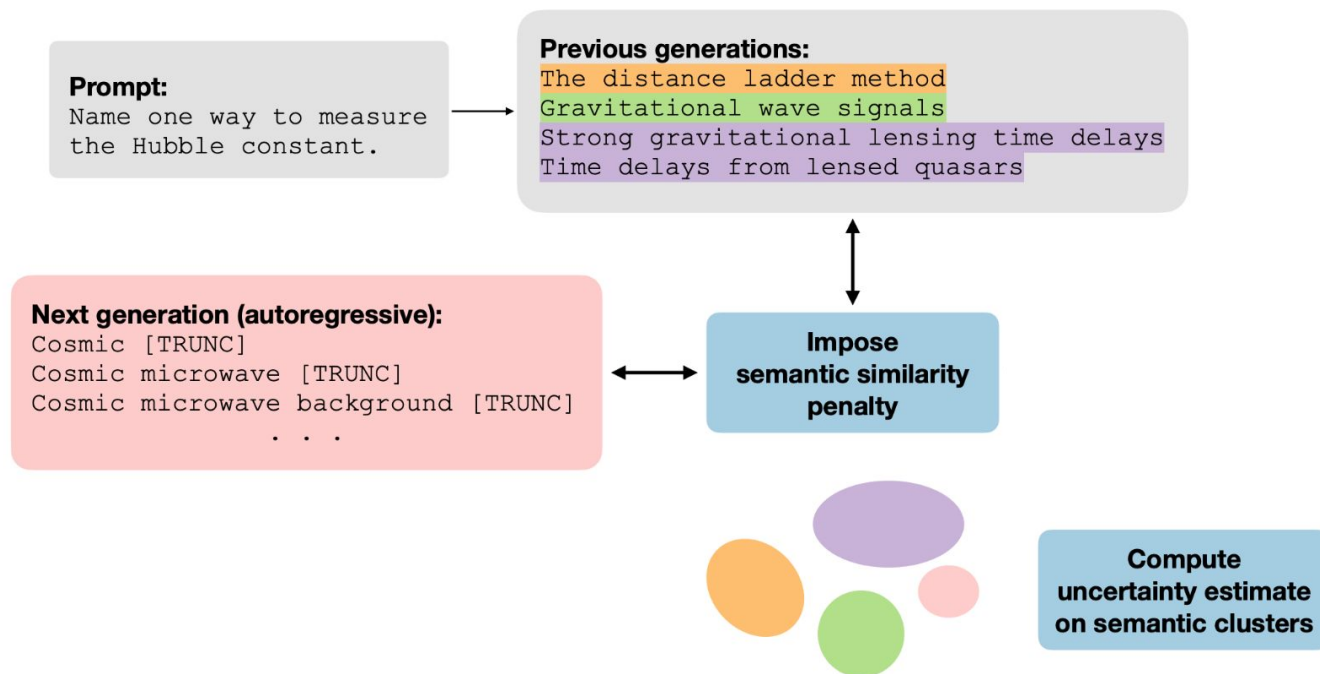
A: The Hubble constant can be measured using the light curves of Type Ia **supernovae**, which have a constant maximum brightness allowing for distance estimation.

A: The method of standard candles involves using **supernovae** as “candles” to measure distances and then infer the Hubble constant from redshift data.

A: One way to measure the Hubble constant is by using the distance-luminosity relationship for Type Ia **supernovae**.

Solution: expose diverse semantic clusters

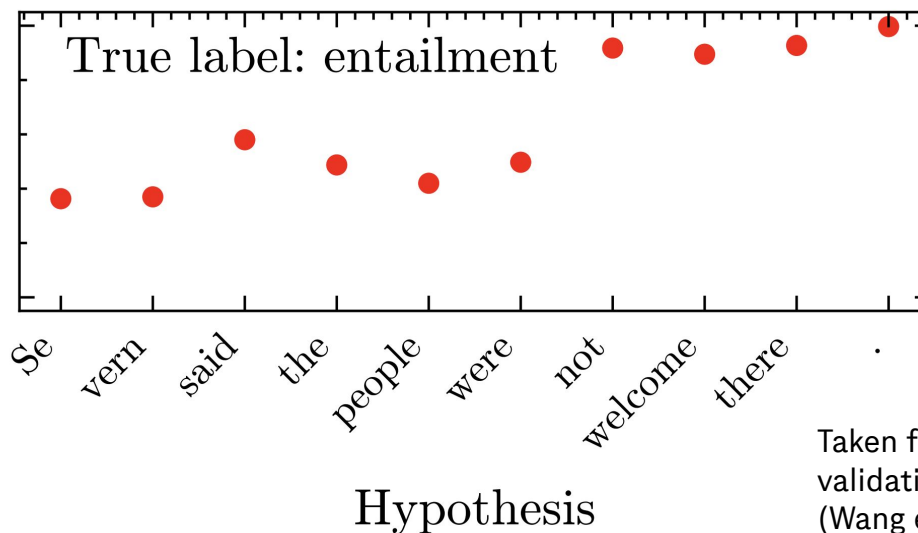
We propose a **diversity-steered sampler** that actively discourages outputs that are semantically redundant to previous generations.



The main ingredient

We introduce a new [TRUNC] token and lightly **finetune** a natural language inference (NLI) model on **partial prefixes** of either the premise or hypothesis at all levels of truncation.

Premise:
You and your friends are not welcome here, said Severn.



Taken from the GLUE MNLI
validation_matched split
(Wang et al. 2018)

Key idea: tilting the token-level conditional distributions

The **bidirectional entailment score** predicted by the finetuned NLI model is used to tilt the sampling distribution away from previous generations.

partial generation and existing sample s assumed to be semantically equivalent if they entail each other

$$E(y_{\leq t}, s) = 1/2 (\text{entailment}(y_{\leq t}, s) + \text{entailment}(s, y_{\leq t}))$$

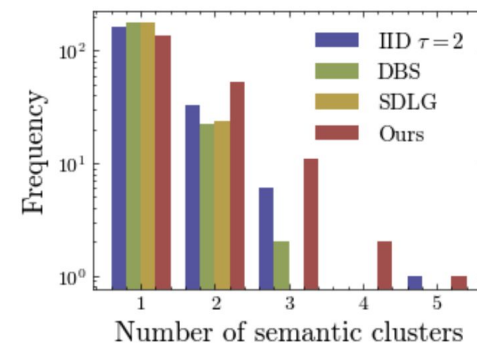
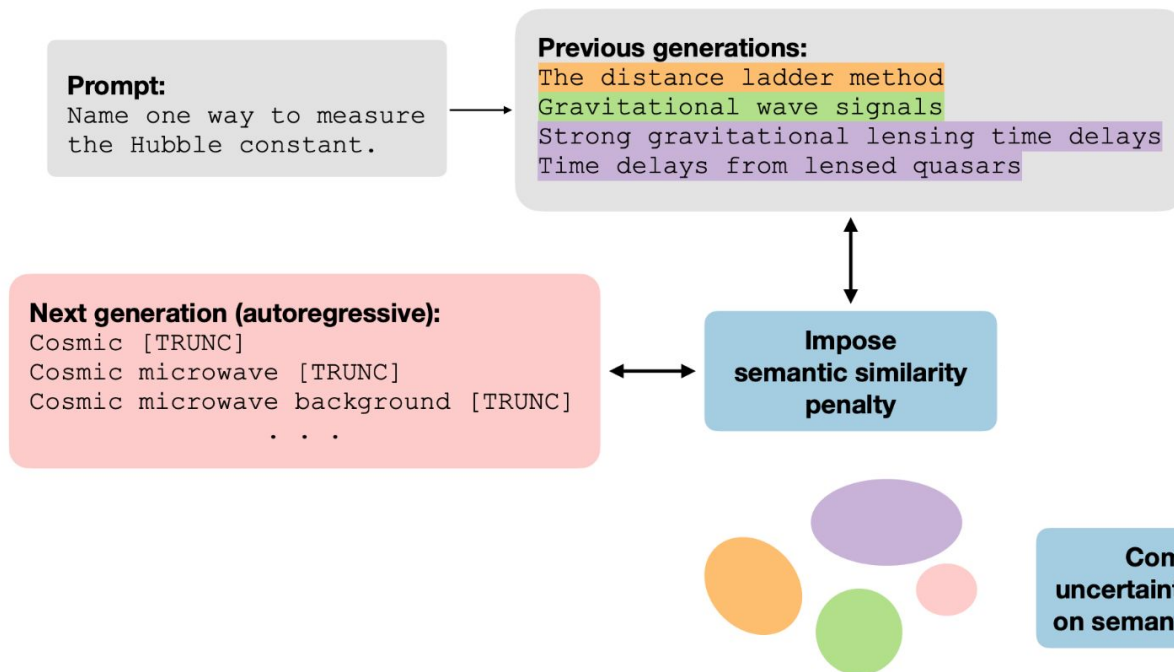
$$\log p(y_t | y_{< t}) - \lambda \max_{s \in \mathcal{S}} E(y_{\leq t}, s)$$

original distribution for token t

repel sampling away from the most similar existing generation

Importance weighting

Sampling from the tilted distribution exposes more semantic clusters with the same number of generations. For estimating semantic entropy downstream, we correct the induced sampling bias with self-normalized importance weights.



$$w_i = \frac{p(s^{(i)})}{q(s^{(i)})} \quad \begin{array}{l} \text{original} \\ \text{tilted} \end{array}$$

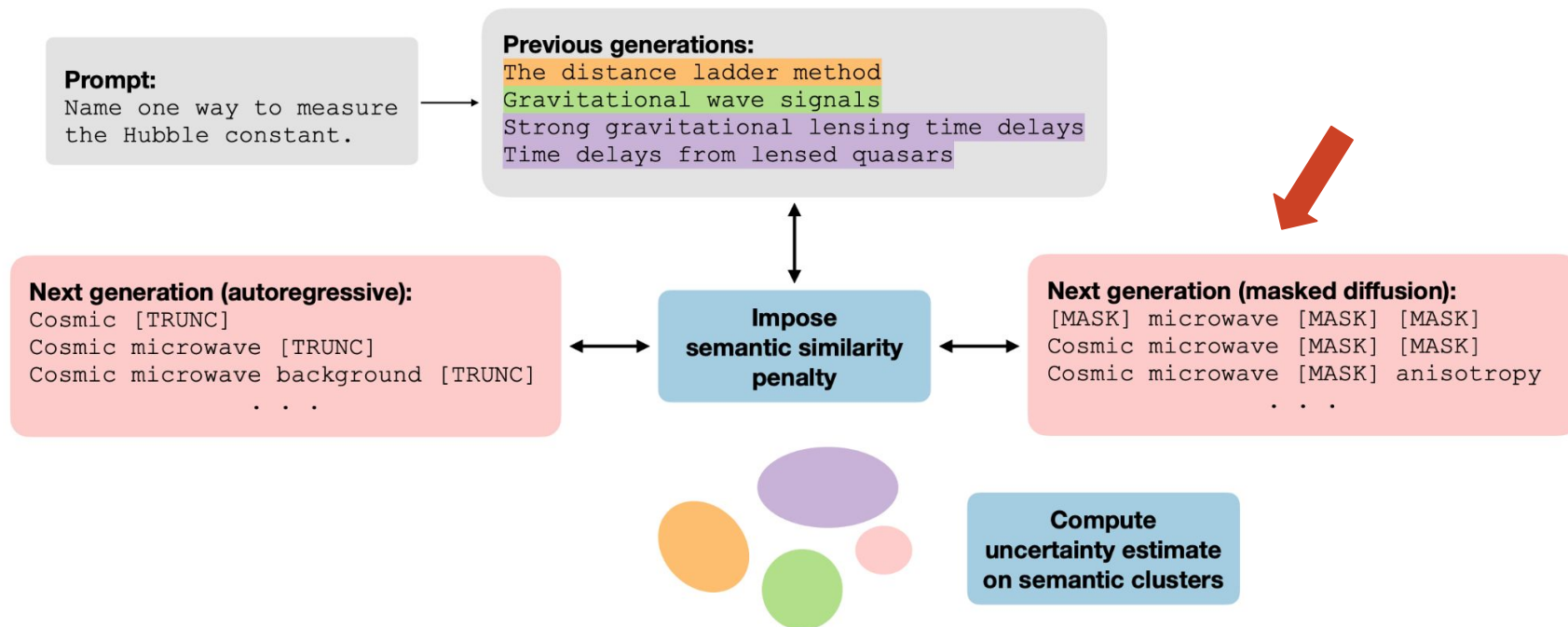
$$\tilde{w}_i = \frac{w_i}{\sum_{j=1}^N w_j}$$

Across four QA benchmarks, our method achieves comparable or better AUROC relative to baselines, including temperature scaling, diverse beam search (Vijayakumar et al., 2018), and SDLG (Aichberger et al., 2025).

Dataset	Model	Vanilla ($\tau = 1$)	$\tau = 2$	DBS [39]	SDLG [4]	Ours
CoQA	OPT-6.7B	.59 \pm .06	.69 \pm .04	.68 \pm .04	.71 \pm .02	.75\pm.02
	OPT-13B	.70 \pm .04	.76\pm.04	.73 \pm .04	.72 \pm .02	.75 \pm .03
	LLaMA 3 8B-Instruct	.68 \pm .03	.72 \pm .04	.71 \pm .05	.74 \pm .02	.77\pm.02
	LLaDA 8B-Instruct	.78 \pm .02	.81\pm.05	-	-	.81\pm.04
TriviaQA	OPT-6.7B	.66 \pm .05	.67 \pm .06	.71 \pm .04	.78 \pm .03	.82\pm.03
	OPT-13B	.72 \pm .04	.70 \pm .05	.73 \pm .04	.86\pm.03	.85 \pm .03
	LLaMA 3 8B-Instruct	.79 \pm .04	.70 \pm .04	.70 \pm .03	.79 \pm .04	.84\pm.03
	LLaDA 8B-Instruct	.81 \pm .11	.83 \pm .05	-	-	.86\pm.04
AmbigQA	OPT-13B	.65 \pm .10	.68 \pm .11	.78\pm.08	.71 \pm .08	.78\pm.04
	LLaMA 3 8B-Instruct	.70 \pm .04	.55 \pm .07	.71 \pm .08	.77\pm.05	.76 \pm .03
	LLaDA 8B-Instruct	.70 \pm .09	.71 \pm .08	-	-	.76\pm.03
TruthfulQA	OPT-6.7B	.80 \pm .04	.80 \pm .05	.77 \pm .02	.78 \pm .06	.81\pm.06
	OPT-13B	.73 \pm .06	.74 \pm .08	.79 \pm .05	.81 \pm .04	.85\pm.04
	LLaMA 3 8B-Instruct	.88 \pm .04	.88 \pm .05	.89\pm.04	.86 \pm .04	.89\pm.02
	LLaDA 8B-Instruct	.85 \pm .04	.89 \pm .04	-	-	.94\pm.02

Both autoregressive and masked diffusion paradigms

Our method also applies to **masked diffusion** models, by similarly finetuning the NLI model on intermediate diffusion states (masked spans).



- We propose an importance sampling framework that actively promotes semantic novelty with respect to previous generations during decoding.
- The key idea is to inject a token-level semantic similarity penalty into the proposal distribution. An NLI model, finetuned minimally with a new [TRUNC] token for autoregressive models, enables live semantic scoring.
- For downstream uncertainty estimation, the induced sampling bias is corrected with self-normalized importance weighting.
- Across four QA benchmarks, our method matches or surpasses baselines while covering more semantic clusters with the same number of samples.
- Our method is modular, requires no gradient access to the base LLM, and applies to both autoregressive and masked diffusion models.

Thank you!

Aichberger, Lukas, et al. "Improving uncertainty estimation through semantically diverse language generation." *ICLR* (2025).

Kuhn, Lorenz, Yarin Gal, and Sebastian Farquhar. "Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation." *ICLR* (2023).

Vijayakumar, Ashwin K., et al. "Diverse beam search: Decoding diverse solutions from neural sequence models." *AAAI* (2018).

Wang, Alex, et al. "GLUE: A multi-task benchmark and analysis platform for natural language understanding." *EMNLP workshop BlackboxNLP* (2018).