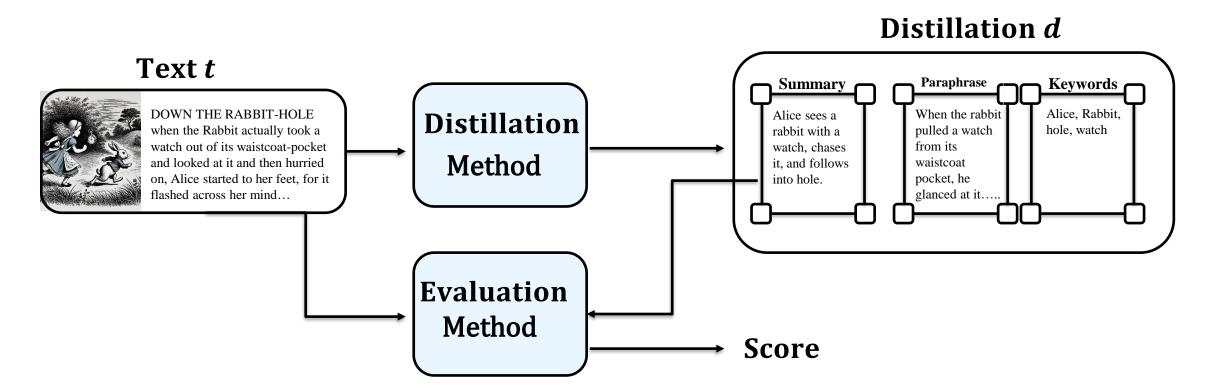
Guided Word Guessing with Language Models and Information-Theoretic Content Distillation



Dana Levin Alon Kipnis

1. Overview

• Content distillation challenges: text summarization, paraphrasing, keyword extraction,...



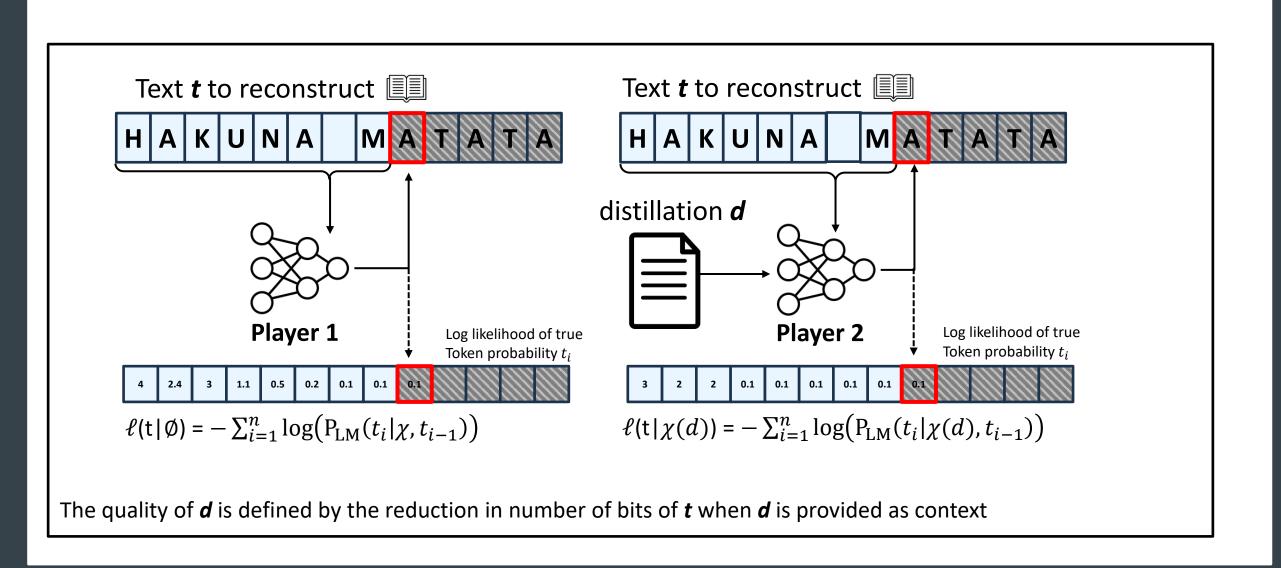
• The **Log Likelihood Gain (LLG)** is the log of the likelihood ratio of **t** under a language model without context to with context **d**:

$$LLG(t; d) := \log[\frac{P_{LM}(t|\chi(d))}{P_{LM}(t|\emptyset)}]$$

• Tenendum: LLG is a simple, natural and useful way to measure the relevance of the distilled content d to the source text t.

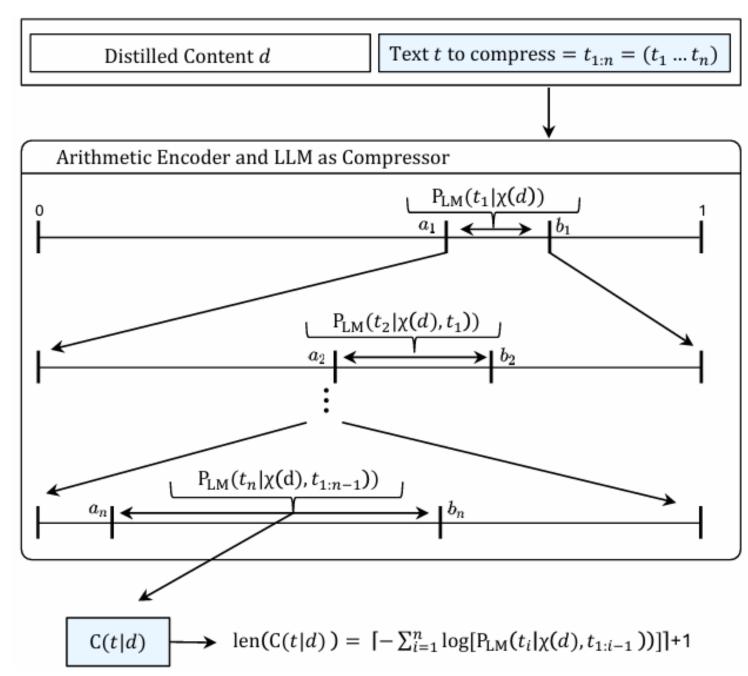
2. Motivation: LLG as a Guessing Game 🞮

- **LLG** is inspired by a variation of the Shannon Guessing game (1951) where a text's information is measured by how well a **human** player predicts it. (Shannon 1951, Hovy et.al. 1998)
- LLG plays the Guessing Game with LMs as players.



3. Lossless Compression

• Result 1: LLG is the reduction in binary codelength when compressing text \boldsymbol{t} using a state-of-the-art lossless compression procedure and distillation \boldsymbol{d} as side information



 $|\operatorname{Len}(C(t|\emptyset)) - \operatorname{Len}(C(t|\chi(d))) - \operatorname{LLG}(t;d)| \leq 1$

(Izacard 2019), (Bellard 2019), (Goyal et. al. 2021), (Mao et. al. 2022), (Deletang et. al. 2023)

4. LLG as Normalized Compression Distance

• NCD (Vitányi 2008) is a universal compression based similarity measure between two objects from any domain:

$$NCD_{Z}(x,y) = \frac{Z(xy) - \min\{Z(x),Z(y)\}}{\max\{Z(x),Z(y)\}}$$

(Cebrián .et.al 2007), (Cilibras, Vitányi 2004), (Pinho .et.al 2016), (Jiang Z .etl.al 2023)

Result 2:

 $\frac{\text{LLG}(t;d)}{\log[1/P_{LM}(t|\emptyset)]}$ is an NCD (under LM + arithmetic encoder compression)

5. Empirical Analysis

Result 3: LLG is empirically competitive with other evaluation methods in the domains of summary and paraphrase evaluation.

Measure	summEval	Newsroom	SFF	CNN	EAPoe	HPotter	CNN-para	HPotter-para
LLG	0.18 (0.017)	0.54 (0.026)	0.26 (0.011)	0.61 (0.01)	0.55 (0.02)	0.6 (0.02)	0.4 (0.03)	0.32 (0.067)
$\overline{ ext{LLG}}$	0.15 (0.018)	0.45 (0.027)	0.25 (0.011)	0.55 (0.012)	0.55 (0.019)	0.61 (0.02)	0.68 (0.008)	0.34 (0.074)
NCD_gzip	0.18 (0.017)	0.44 (0.027)	0.26 (0.011)	0.45 (0.013)	0.48 (0.023)	0.56 (0.024)	0.37 (0.032)	-0.45 (0.065)
BLANC	0.13 (0.018)	0.48 (0.026)	0.26 (0.011)	0.39 (0.014)	0.3 (0.027)	0.36 (0.031)	0.2 (0.037)	-0.31 (0.073)
BARTScore	0.27 (0.016)	0.5 (0.028)	0.15 (0.012)	0.43 (0.014)	0.06 (0.031)	0.06 (0.035)	-0.71 (0.001)	-0.71 (0.005)
ROUGUEL	0.13 (0.017)	0.49 (0.027)	0.27 (0.011)	0.48 (0.012)	0.46 (0.022)	0.6 (0.02)	0.71 (0.002)	0.42 (0.065)
BERTScore	0.31 (0.016)	0.48 (0.03)	0.2 (0.011)	0.63 (0.01)	0.56 (0.019)	0.61 (0.021)	-0.71 (0.001)	-0.71 (0.005)

(T Zhang .et.al 2019), (CY Lin 2004), (W Yuan .et.al 2021), (O Vasilyev .etl.al 2020)

6. Structure/Content Word Analysis

• **Result 4:** LLG is overwhelmingly affected by the text's content words than its structure – desired behavior for distillation evaluation.

					0.0100 -				C		tokens e tokens	
Measure	summEval	Newsroom	SFF	CNN	e 0.00/3		l				- 10/10/10	
LLG	0.18 (0.018)	0.54 (0.026)	0.26 (0.011)	0.61 (0.01)	ਲੂ 0.0050 -		الهاا	4				
LLG _{cont}	0.17 (0.017)	0.54 (0.026)	0.27 (0.01)	0.6 (0.011)	0.0025 -							
LLG _{stru}	0.15 (0.017)	0.47 (0.03)	0.15 (0.011)	0.34 (0.015)	0.0025					L		
Kendall tau correlation with reference scores					0.0000	Ō	100	200 LLC	300 [bits]	400	500	I

7. Content Distillation with LLG

• Can we use LLG as an objective for content distillation?

$$d^*(t) = \underset{d \in D}{\operatorname{argmax}} \operatorname{LLG}(d; t)$$

To find $d^*(t)$, we measure LLG for each potential distillation $d\epsilon D$

Slogan: <u>d</u>. Q: What is the product? A: Organic skin-care line.

Slogan_d. Q: What are the brand values? A: Sustainability, purity, self care.

Slogan <u>d</u>. Q: Who is the target audience? A: Environmentally conscious individuals seeking natural beauty solutions.

Slogan___d__. Q: What is the brand tone? A: Gentle and nurturing.

For $D = \{\text{all word bigram}\}, d^*(t) = \text{Naturally Blossom}\}$

• $d^*(t)$ is <u>objectively</u> the optimal distillation from the LM perspective.

Challenge		LLG-optimized distilled content				
Tagline	Ever noticed ho of people taking putting passeng	Overcrowd Aircraft				
Tagline	A drunk teenage boy had to be rescued by security after jumping into a lions' enclosure at a zoo in western India.					
	What is the	What are the	Who is the target audience?	What is the brand		
	product?	brand values?		tone?		
Slogan	Organic skin- care line	Sustainability, pu- rity, self care	Environmentally conscious indi- viduals seeking natural beauty solutions	Gentle and nurturing	Naturally Blossom	
Slogan	Luxury watch	Craftsmanship, tradition, prestige	High income individuals who value timeless elegance	Classic and refined	timeless elegance	
Slogan	Energy drink.	Performance, vi- tality, boldness	Active individuals and fitness enthusiasts.	Energetic and high spirited	Energetic high	
Slogan	Running shoes.	Performance, durability, innova- tion	Runners and athletes who aim to push their limits	Motivational and en- ergetic	Nike mph	
Slogan	Streaming service	Convenience, en- tertainment, vari- ety	People seeking on-demand video content	Engaging and inclusive	Hulu on-demand	
Slogan	Electric car	Sustainability, in- novation, luxury	Environmentally conscious drivers who value performance.	Sophisticated and forward thinking	Tesla S	