

Estimating Text Temperature

Anonymous ACL submission

Abstract

Autoregressive language models typically use temperature parameter at inference to shape the probability distribution and control the randomness of the text generated. After the text was generated, this parameter can be estimated using maximum likelihood approach. Following it, we propose a procedure to estimate the temperature of any text, including ones written by humans, with respect to a given language model. We evaluate the temperature estimation capability of a wide selection of small-to-medium LLMs. We then use the best-performing Qwen3 14B to estimate temperatures of popular corpora.

1 Text Generation with a Language Model

Given an input text as a context, the goal of open-ended generation is to produce a coherent continuation of the text (Holtzman et al., 2020). More formally, given a sequence of m tokens $t^{(1)} \dots t^{(m)}$ as a context, the objective is to generate the next n continuation tokens, resulting in the completed sequence $t^{(1)} \dots t^{(m+n)}$. This is achieved through the use of the left-to-right text probability decomposition, which is used to generate the sequence one token at a time, using a particular decoding strategy.

A common approach to text generation is to shape a probability distribution through temperature (Ackley et al., 1985). Let

- $u^{(i)}$ be the vector of logits at the step i .
- $u_{obs}^{(i)}$ be the specific logit value corresponding to the token $t^{(i)}$.

The probability of observing the specific token at step i is then re-estimated as:

$$p(t^{(i)}|u^{(i)}, T) = \frac{\exp(u_{obs}^{(i)}/T)}{\sum_l \exp(u_l^{(i)}/T)} \quad (1)$$

Setting $T \in [0, 1)$ skews the distribution towards high-probability events, and, similarly, $T \in (1, \infty)$ skews the distribution towards low-probability events.

2 Estimating Temperatures

2.1 Estimating the Temperature of a Generated Sequence

Suppose that we have logits for multiple generation steps and actual generated tokens for the same steps. What is the maximum likelihood estimate of temperature T ? We can treat it as a classic parameter estimation problem.

The total log-likelihood $\mathcal{L}(T)$ for all N steps is the sum of the individual log-probabilities:

$$\mathcal{L}(T) = \sum_{i=1}^N \log \left(\frac{\exp(u_{obs}^{(i)}/T)}{\sum_l \exp(u_l^{(i)}/T)} \right) \quad (2)$$

To find the MLE, we take the derivative of $\mathcal{L}(T)$ with respect to T and set it to zero.

$$\begin{aligned} \frac{d\mathcal{L}}{dT} &= \\ \sum_{i=1}^N \left[-\frac{u_{obs}^{(i)}}{T^2} - \left(-\frac{1}{T^2} \sum_l u_l^{(i)} p(l|u^{(i)}, T) \right) \right] & \\ &= \frac{1}{T^2} \sum_{i=1}^N \left[\mathbb{E}[u^{(i)}|T] - u_{obs}^{(i)} \right] \end{aligned} \quad (3)$$

Setting the derivative (3) to zero, we get for the MLE \hat{T} :

$$\sum_{i=1}^N u_{obs}^{(i)} = \sum_{i=1}^N \mathbb{E}[u^{(i)}|\hat{T}], \quad (4)$$

where

$$\mathbb{E}[u^{(i)}|\hat{T}] = \sum_l u_l^{(i)} \frac{\exp(u_l^{(i)}/\hat{T})}{\sum_k \exp(u_k^{(i)}/\hat{T})} \quad (5)$$

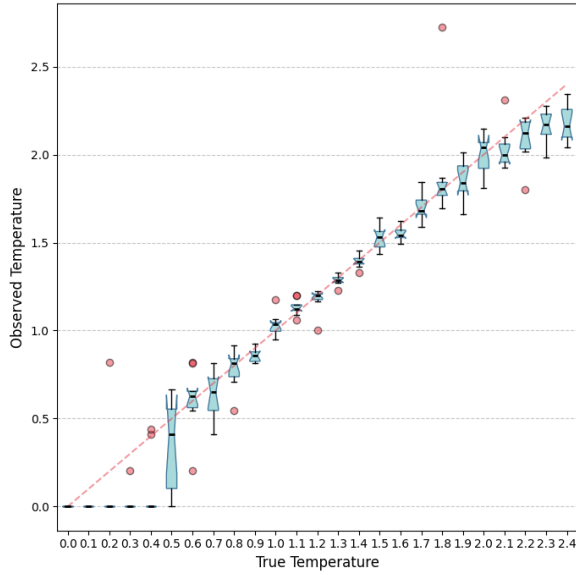


Figure 1: Estimated vs. generation temperature for granite-4.0-micro.

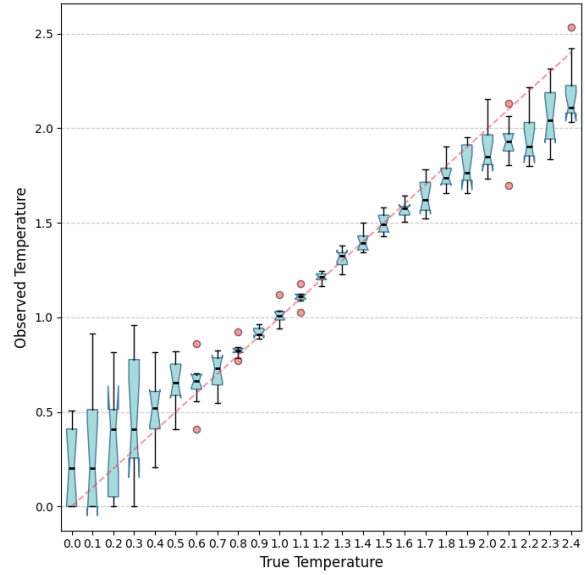


Figure 2: Estimated vs. generation temperature for Llama-3.2-3B.

Thus, the maximum likelihood estimate for T is the temperature at which the sum of the observed logits equals the sum of the expected logits predicted by the model.

2.2 Estimating the Temperature of Any Text

Now, observe that the text we work on in the previous subsection does not need to be a generated one. As soon as we have tokenized a text and have estimated the probability distributions for all the tokens in the text with some probabilistic model (for example, an LLM), we can estimate the text temperature using (4). This estimate will, of course, depend on the LLM model used.

To specifically find the temperature, we numerically solve equation (4) for T , given the logits for all the tokens and specific tokens present in the text.

3 Estimating Temperature of LLM-Generated Texts

To better understand the properties of the proposed temperature estimation approach, we generate texts with an assortment of LLMs in small to medium sizes representing different LLM families, finetuning levels and architectures:

- Qwen3 0.6B, 1.7B 4B, 8B, 14B (Yang et al., 2025) (normal, base and 4B-Thinking-2507)
- Phi-4-mini (Microsoft Team, 2025) (reasoning and instruct)

- Llama 3.1 8B, 3.2 1B, 3B, 11B-Vision (LLama Team, 2024) (base and instruct)
- gemma-3 270m, 1b, 4b, 12b (Gemma Team, 2025)
- DeepSeek-R1 distills: 0528-Qwen3-8B, Qwen-1.5B, Qwen-7B, Qwen-14B, and Llama-8B (DeepSeek AI, 2025)
- granite-4.0-micro (IBM Research, 2025) (base and normal)

We use 4-bit BitsAndBytes NF4 quantized models through the HuggingFace Transformers library (Wolf et al., 2020). All the texts are generated from a single random seed-controlled token using batch inference. We generate 10 texts 200 tokens long for each temperature in a range from 0.001 to 2.401 with step 0.1. We do not use top-k, top-p or any other decoding parameters such as no-repeat. To estimate the temperature, we use SciPy (Virtanen et al., 2020) root finding function `root_scalar` for the reverse temperature with bracket $[10^{-2}, 10^4]$ and initial value $5 * 10^3$. To speed up the computations we compute softmax (1, 5) using PyTorch on GPU. The code was mostly written by Gemini.

3.1 What Does the Measured Temperature Say About the Generation Temperature?

In the first series of experiments, we estimate the temperature of texts generated by an LLM with the

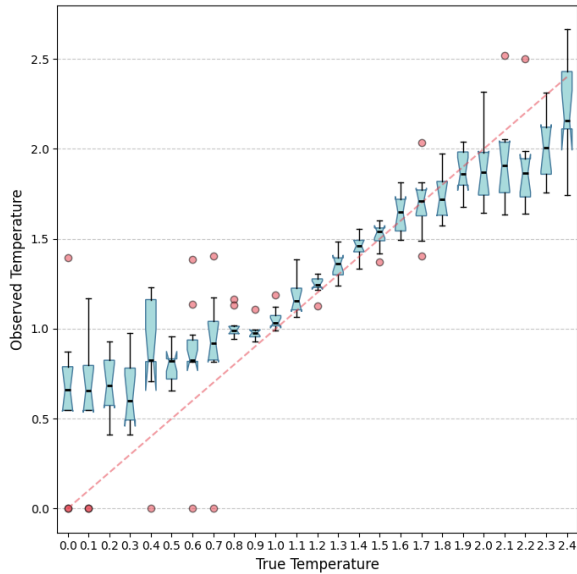


Figure 3: Qwen3-8B estimates Qwen3-4B

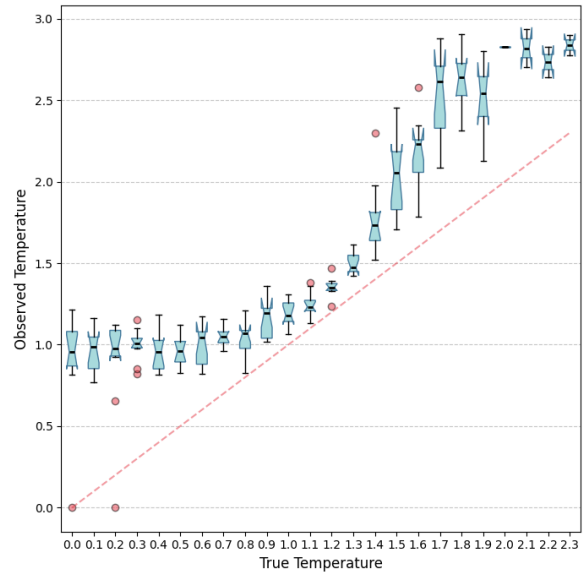


Figure 4: DeepSeek-R1-0528-Qwen3-8B estimates DeepSeek-R1-Distill-Qwen-1.5B

116 same LLM. For most LLMs, the estimated temper-
 117 ature insignificantly differs from the generation
 118 temperature.

119 For all the models, we observe pronounced and
 120 systematic differences between estimated and gen-
 121 eration temperatures in areas of low and high tem-
 122 peratures (see Figures 1, 2). The difference in the
 123 low temperature area is the result of saturation
 124 when the generation with different temperatures
 125 produces one and the same result. Thus, the tem-
 126 perature estimation at low temperatures is an ill-
 127 posed problem that can be resolved using some kind
 128 of regularization yet to find. This constitutes an excit-
 129 ing area for future research. We don't yet have an
 130 explanation for the discrepancy between estimated
 131 and generation temperatures at high temperatures.

132 3.2 What Does the Measured Temperature 133 say About the Model?

134 In the second series of experiments, we estimate
 135 the temperature with an LLM different from the
 136 one used to generate the text. The observed behav-
 137 iors can be put into several qualitatively distinct
 138 groups. In one group, the graph of estimated vs.
 139 true temperature does not differ much from what a
 140 graph would be for the same LLM. This happens
 141 mostly for similar models in the families of Qwen,
 142 granite and LLama (see Figure 3 for an example).

143 The other group consistently overestimates the
 144 generation temperature (see Figure 4 for an ex-
 145 ample). This is typical for DeepSeek distills and
 146 happens in some other cases, but never for base

models. We attribute this to narrower probabil-
 ity distributions of finetuned, especially reasoning,
 models.

Yet another group consists of model pairs that
 produce a pronounced S-shaped graph (see Figure
 5). Finally, there are model pairs that show
 little to no correlation between the generation and
 the estimated temperatures (see Figure 6 for an ex-
 ample). This is often the case when either generator
 or estimator models are gemma or Phi.

To digitize these observations, we calculated sta-
 tistical metrics of goodness such as MAE, R^2 , Pear-
 son p between observed and generation tempera-
 tures for each pair of models. Figure 7 presents a
 heatmap of MAE between the models. From this
 heatmap we can conclude that:

- Models in Qwen, LLama and granite fami-
 lies reasonably well estimate temperature of
 models from these families. They estimate
 even better when the generator and estimator
 models are from the same family.
- DeepSeek, gemma and Phi models are not
 good at both estimating and being estimated
- Base models are overall slightly better tem-
 perature estimators than finetuned ones. Base
 models better estimate base models and fine-
 tuned better estimate finetuned ones. This is
 especially pronounced for LLama family.
- Larger models are slightly better estimators

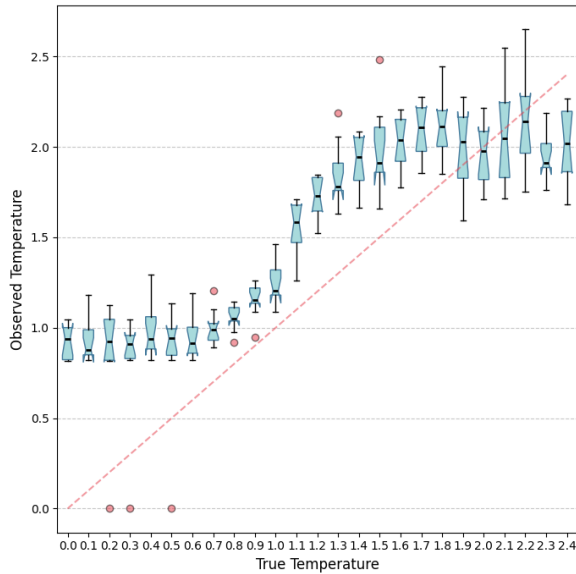


Figure 5: DeepSeek-R1-0528-Qwen3-8B estimates Meta-Llama-3.1-8B-Instruct

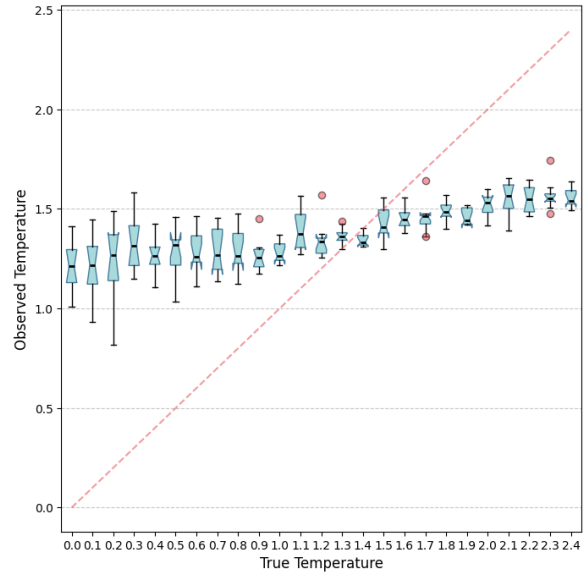


Figure 6: gemma-3-1b-it estimates Qwen3-14B

Dataset	Mean t	std
WikiText	1.0585	0.0645
Poetry	1.0089	0.1044
Jokes	1.1003	0.1037
GSM8K	1.0942	0.0992
Code/Python	0.9242	0.1539
IMDB	1.0244	0.0385
HH RLHF	1.0063	0.0641
AG News	1.1023	0.0734
Yelp	1.0349	0.0545

Table 1: Temperatures of corpora

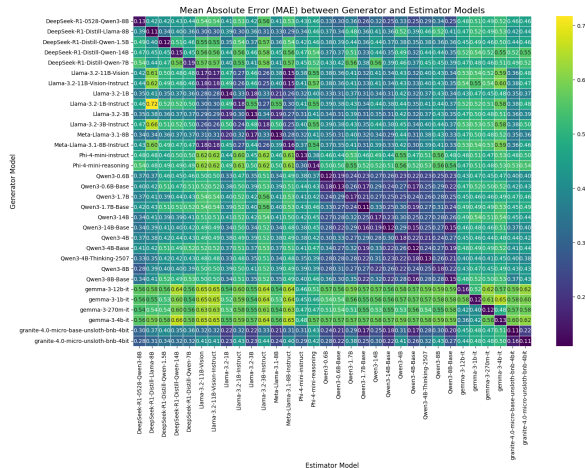


Figure 7: Heatmap of MAE between models

- Qwen models are the best estimators, the absolute leader being Qwen3-14B-base

4 Estimating Temperature of Human-written Corpora

Based on the results in the previous section, we use Qwen3 14B Base as the preferred LLM for the temperature estimation. We estimate temperature of several diverse datasets: WikiText-103 (Merity et al., 2016), Poetry dataset (Unknown, 2025b), Jokes dataset (Unknown, 2025a), GSM8K (Cobbe et al., 2021), Python code (Unknown, 2025c), IMDB (Maas et al., 2011), HH RLHF (Bai et al., 2022; Ganguli et al., 2022), AG News (Zhang et al., 2015b; Gulli, 2005; del Corso et al., 2005), Yelp (Zhang et al., 2015a). We sample 300 texts from each dataset and average temperatures over texts in a dataset. The results are in the Table 1.

The estimated temperatures are all close to 1. This is unsurprising. Notable exceptions are Jokes, GSM8K and AG News that have a statistically significant higher temperature of 1.1, and Python code that has a lower temperature of 0.9. Explaining these temperature phenomena is an exciting topic for future research.

5 Conclusion

We proposed a procedure to estimate the temperature of any text, with respect to a given language model. We evaluated the temperature estimation capability of a wide selection of small-to-medium LLMs. While the temperatures measured for most corpora are close to 1, the suggested method spots some interesting temperature phenomena.

208	Limitations		
209	The experiments were conducted for English only.		
210	This may limit generalization to other languages.		
211	The list of LLM tested is limited to small-to-		
212	medium models. Scaling laws, if any, are unclear.		
213	References		
214	David H. Ackley, Geoffrey E. Hinton, and Terrence J.		
215	Sejnowski. 1985. A learning algorithm for boltz-		
216	mann machines . <i>Cognitive Science</i> , 9(1):147–169.		
217	Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda		
218	Askell, Anna Chen, Nova DasSarma, Dawn Drain,		
219	Stanislav Fort, Deep Ganguli, Tom Henighan,		
220	Nicholas Joseph, Saurav Kadavath, Jackson Kernion,		
221	Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac		
222	Hatfield-Dodds, Danny Hernandez, Tristan Hume,		
223	Scott Johnston, Shauna Kravec, Liane Lovitt, Neel		
224	Nanda, Catherine Olsson, Dario Amodei, Tom		
225	Brown, Jack Clark, Sam McCandlish, Chris Olah,		
226	Ben Mann, and Jared Kaplan. 2022. Training		
227	a helpful and harmless assistant with reinforce-		
228	ment learning from human feedback . <i>Preprint</i> ,		
229	arXiv:2204.05862.		
230	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,		
231	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias		
232	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro		
233	Nakano, Christopher Hesse, and John Schulman.		
234	2021. Training verifiers to solve math word prob-		
235	lems. <i>arXiv preprint arXiv:2110.14168</i> .		
236	DeepSeek DeepSeek AI. 2025. Deepseek-r1: Incen-		
237	tivating reasoning capability in llms via reinforce-		
238	ment learning . <i>Preprint</i> , arXiv:2501.12948.		
239	Gianna Maria del Corso, Antonio Gulli, and Francesco		
240	Romani. 2005. Ranking a stream of news . In <i>The</i>		
241	<i>Web Conference</i> .		
242	Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda		
243	Askell, Yuntao Bai, Saurav Kadavath, Ben Mann,		
244	Ethan Perez, Nicholas Schiefer, Kamal Ndousse,		
245	Andy Jones, Sam Bowman, Anna Chen, Tom Con-		
246	erly, Nova DasSarma, Dawn Drain, Nelson Elhage,		
247	Sheer El-Showk, Stanislav Fort, Zac Hatfield-Dodds,		
248	Tom Henighan, Danny Hernandez, Tristan Hume,		
249	Josh Jacobson, Scott Johnston, Shauna Kravec,		
250	Catherine Olsson, Sam Ringer, Eli Tran-Johnson,		
251	Dario Amodei, Tom Brown, Nicholas Joseph, Sam		
252	McCandlish, Chris Olah, Jared Kaplan, and Jack		
253	Clark. 2022. Red teaming language models to re-		
254	duce harms: Methods, scaling behaviors, and lessons		
255	learned . <i>Preprint</i> , arXiv:2209.07858.		
256	Google Gemma Team. 2025. Gemma 3 technical report .		
257	<i>Preprint</i> , arXiv:2503.19786.		
258	A. Gulli. 2005. The anatomy of a news search engine .		
259	In <i>Special Interest Tracks and Posters of the 14th</i>		
260	<i>International Conference on World Wide Web</i> , WWW		
261	'05, page 880–881, New York, NY, USA. Association		
262	for Computing Machinery.		
	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and	263	
	Yejin Choi. 2020. The curious case of neural text de-	264	
	generation . In <i>International Conference on Learning</i>	265	
	<i>Representations</i> .	266	
	IBM Research. 2025. Granite 4.0 language	267	
	models. https://github.com/ibm-granite/	268	
	granite-4.0-language-models . Accessed: 2025-	269	
	10-01.	270	
	Meta Llama Team. 2024. The llama 3 herd of models .	271	
	<i>Preprint</i> , arXiv:2407.21783.	272	
	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,	273	
	Dan Huang, Andrew Y. Ng, and Christopher Potts.	274	
	2011. Learning word vectors for sentiment analysis .	275	
	In <i>Proceedings of the 49th Annual Meeting of the</i>	276	
	<i>Association for Computational Linguistics: Human</i>	277	
	<i>Language Technologies</i> , pages 142–150, Portland,	278	
	Oregon, USA. Association for Computational Lin-	279	
	guistics.	280	
	Stephen Merity, Caiming Xiong, James Bradbury, and	281	
	Richard Socher. 2016. Pointer sentinel mixture mod-	282	
	els . <i>Preprint</i> , arXiv:1609.07843.	283	
	Microsoft Microsoft Team. 2025. Phi-4-mini tech-	284	
	nical report: Compact yet powerful multimodal	285	
	language models via mixture-of-loras . <i>Preprint</i> ,	286	
	arXiv:2503.01743.	287	
	Unknown. 2025a. Jokes dataset .	288	
	Unknown. 2025b. Poetry dataset .	289	
	Unknown. 2025c. python code dataset .	290	
	Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt	291	
	Haberland, Tyler Reddy, David Cournapeau, Ev-	292	
	geni Burovski, Pearu Peterson, Warren Weckesser,	293	
	Jonathan Bright, Stéfan J. van der Walt, Matthew	294	
	Brett, Joshua Wilson, K. Jarrod Millman, Nikolay	295	
	Mayorov, Andrew R. J. Nelson, Eric Jones, Robert	296	
	Kern, Eric Larson, C J Carey, Ilhan Polat, Yu Feng,	297	
	Eric W. Moore, Jake VanderPlas, Denis Laxalde,	298	
	Josef Perktold, Robert Cimrman, Ian Henriksen, E. A.	299	
	Quintero, Charles R. Harris, Anne M. Archibald, An-	300	
	tônio H. Ribeiro, Fabian Pedregosa, Paul van Mul-	301	
	bregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0:	302	
	Fundamental Algorithms for Scientific Computing in	303	
	Python . <i>Nature Methods</i> , 17:261–272.	304	
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	305	
	Chaumond, Clement Delangue, Anthony Moi, Pier-	306	
	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	307	
	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	308	
	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	309	
	Teven Le Scao, Sylvain Gugger, Mariama Drame,	310	
	Quentin Lhoest, and Alexander Rush. 2020. Trans-	311	
	formers: State-of-the-art natural language processing .	312	
	In <i>Proceedings of the 2020 Conference on Empirical</i>	313	
	<i>Methods in Natural Language Processing: System</i>	314	
	<i>Demonstrations</i> , pages 38–45, Online. Association	315	
	for Computational Linguistics.	316	

317 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
318 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
319 Chengen Huang, Chenxu Lv, Chujie Zheng, Dayi-
320 heng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge,
321 Haoran Wei, Huan Lin, Jialong Tang, Jian Yang,
322 Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi
323 Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai
324 Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao
325 Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang,
326 Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan
327 Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao
328 Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xu-
329 ancheng Ren, Yang Fan, Yang Su, Yichang Zhang,
330 Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang,
331 Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zi-
332 han Qiu. 2025. [Qwen3 technical report](#). *Preprint*,
333 arXiv:2505.09388.

334 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015a.
335 Character-level convolutional networks for text clas-
336 sification. In *Proceedings of the 29th International*
337 *Conference on Neural Information Processing Sys-*
338 *tems - Volume 1*, NIPS'15, page 649–657, Cambridge,
339 MA, USA. MIT Press.

340 Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015b.
341 Character-level convolutional networks for text clas-
342 sification. In *NIPS*.