TracVC: Tracing Verbalized Confidence of LLMs Back to Training Data

Anonymous ACL submission

Abstract

Large language models (LLMs) can increase users' perceived trust by verbalizing confidence in their outputs. However, prior work shows that LLMs often express overconfidence, which is misaligned with factual accuracy. To better understand the sources of this behavior, we propose **TracVC**, a method for **Trac**ing Verbalized Confidence back to specific training data. We conduct experiments on OLMo models in a question answering setting, defining a model as truthful when content-related training data-relevant to the question and answer, has greater influence than confidence-related data. Our analysis reveals that OLMo2-13B is often influenced by confidence-related data that is semantically unrelated to the query, suggesting that it may mimic linguistic markers of certainty. This finding highlights a fundamental limitation in current training regimes: LLMs may learn how to sound confident without understanding when confidence is warranted. Our analysis provides a foundation for improving LLMs' trustworthiness in expressing more truthful confidence.¹

1 Introduction

005

011

015

017

019

021

026

027

034

Verbalized confidence in large language models (LLMs) is increasingly used to estimate the certainty of their generated outputs to improve transparency and user trust (Kadavath et al., 2023; Tian et al., 2023; Chen and Mueller, 2024). However, recent studies have shown that LLMs frequently exhibit overconfidence, which is not aligned with factual accuracy (shown in Fig. 1), resulting in their poor reliability in expressed confidence (Zhou et al., 2024; Xiong et al., 2024; Ni et al., 2025; Xia et al., 2025). These findings lead to a foundational question: what drives confidence in LLMs, and do LLMs understand the intended meaning of expressing confidence?



Figure 1: LLMs tend to be overconfident. E.g., more than 94% test samples get confidence scores between 0.9 to 1.0 from the OLMo2-13B Instruct model, even though the accuracy of those samples is less than 0.4.

041

042

043

045

047

055

058

060

061

062

063

064

065

066

In this paper, we investigate the origin of verbalized confidence in LLMs by examining the influence of training data. Specifically, we ask: Do LLMs ground their confidence on semantically relevant, content-related training samples, or are they instead influenced by superficial confidencerelated cues? To address this question, we propose **TracVC**, a method for tracing the origins of verbalized confidence back to specific training data. TracVC builds on influence estimation techniques (Pruthi et al., 2020), but is adapted to estimate how individual training samples contribute to confidence generation in LLMs.

We apply TracVC to analyze 11 open-source LLMs with publicly accessible training data, OLMo (Team OLMo et al., 2024) and Llama3-8B (Meta AI, 2024) models, on five question answering benchmarks. To measure whether contentrelated training data (i.e., question and answer) is more influential than confidence-related data in shaping verbalized confidence, we define a truthfulness metric as the ratio of cases where content wins over confidence. Our results demonstrate that (1) OLMo2-13B models are influenced more by confidence-related training samples that are semantically unrelated to the question, often latching onto

¹Code is available at: https://anonymous.4open. science/r/training_data_confidence-CB2E.

keywords like "confidence" regardless of context. 067 This suggests that it does not ground its confidence 068 in truthful information, but instead learns to mimic 069 linguistic markers of certainty; (2) Larger LLMs do not exhibit higher truthfulness than smaller ones; (3) LLMs can show low truthfulness even if they may have seen content-related samples during train-073 ing; (4) Post-training techniques (e.g., direct preference optimization) have limited and inconsistent impact on truthfulness, suggesting that pre-training may be more critical in shaping how confidence is grounded.

> These findings highlight a fundamental limitation in current training regimes: LLMs may learn to sound confident without understanding when confidence is warranted. Our work introduces a datadriven perspective of model confidence, offering insights that can guide future training approaches toward improving the truthfulness and reliability of model confidence.

2 Related Work

081

094

098

103

104

105

106

107

108

Overconfidence in LLMs (Zhou et al., 2024; Xiong et al., 2024; Ni et al., 2025; Xia et al., 2025) poses a risk to user trust. Interpreting verbalized confidence is thus crucial to address this risk. Ni et al. (2025); Kumar et al. (2024) only investigates whether verbalized confidence is aligned with internal probabilities of LLMs, lacking insights into the origins of confidence. Influence estimation methods, such as influence functions (Koh and Liang, 2017), offer an alternative to leaveone-out re-training when studying the impact that specific examples in the training data have on a model. This has enabled their recent application to related problems such as fact-tracing (Chang et al., 2024), outlier detection (Pruthi et al., 2020), or data valuation (Choe et al., 2024; Bejan et al., 2023). In this work, we adapt the gradient-similarity based method TracIn (Pruthi et al., 2020) to enable efficient analysis of what training data influences LLM's confidence estimates, an application not previously approached through this lens.

3 Methodology

110The workflow of our proposed TracVC is shown111in Fig. 2. TracVC identifies what types of train-112ing data are more influential for LLM verbalized113confidence. We mainly study two types of training114data: content-related and confidence-related data.115The following sections explain our prompt design,



Figure 2: The workflow of TracVC. We first search the top 10 similar samples for content (question (q) and answer (a)) and confidence (phrase (p) and confidence (c)). Then, we compute and compare the influence score regarding the confidence generation for content-related and confidence-related training data.

the data search mechanism to retrieve related data for content and confidence from training data, and finally, the method we use to estimate and compare the influence score of the retrieved data. 116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

133

134

135

136

137

138

139

140

141

3.1 Prompt Design

We use a two-step prompt for this experiment. The first-stage prompt contains the question q for an LLM \mathcal{M} to generate (\sim) the answer a ($a \sim \mathcal{M}(q)$). The second-stage prompt consists of a phrase p to require \mathcal{M} to provide a probability for its generated answer. Finally, \mathcal{M} generates the verbalized confidence c regarding a and q ($c \sim \mathcal{M}(q, a, p)$).

3.2 Training Data Search

To answer our research question on whether the question-answer (*content*)-related training data is more influential than the *confidence*-related training data, we first search for the top 10 most similar data samples of each type per test sample from the training data. To achieve this, we employ the search mechanism developed by Elazar et al. (2024), which uses Elastic Search (Elastic, 2025) to create a keyword-based search engine built on the training data. This search engine allows us to retrieve similar samples for a search query: By querying with question and answer, we obtain 10 content-related samples $T = \{d_1^t, ..., d_{10}^t\}$. Likewise, by



Figure 3: Examples of the most influential data samples for different LLMs when generating confidence. We retrieve these samples from the pre-training data, we provide additional examples in Appendix A.7.

searching for phrase and confidence, we obtain 10 142 confidence-related samples $F = \{d_1^f, ..., d_{10}^f\}$. 143

3.3 **Training Data Influence Estimation**

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

167

171

172

1

Inspired by TracIn (Pruthi et al., 2020), our method is based on point-wise comparisons of loss gradients. Specifically, we estimate the influence of a training sample d on the model's generated confidence c by computing the cosine similarity between the gradients of the model's loss function. One gradient is computed when the loss is evaluated on the training sample d, and the other when evaluated on (q, a, p). Both gradients are taken with respect to the model's input embeddings w:

$$\operatorname{Inf}(d|c_i) = \frac{\nabla \ell(w_i, d) \cdot \nabla \ell(w_i, (q_i, a_i, p))}{\|\nabla \ell(w_i, d)\| \cdot \|\nabla \ell(w_i, (q_i, a_i, p))\|}$$
(1)

We deliberately exclude the verbalized confidence c from this score to ensure its independence of confidence variations. The implementation details are discussed in Appendix A.2.

3.4 **Truthfulness Evaluation**

We define truthfulness as the property whereby content-related training data exerts a greater influence than confidence-related data on confidence 163 generation. To quantify truthfulness, we intro-164 duce **ccr** (Content-over-Confidence Ratio), defined 165 as the winning counts ratio between the contentrelated set T and confidence-related set F. Intuitively, a higher ccr indicates that the model's confidence is more strongly driven by semantically 170 relevant, truthful content rather than by potentially misleading confidence cues. Assume our test data contains *n* questions, **ccr** is computed by:

73
$$ccr = \frac{\sum_{i=0}^{n} \sum_{d^{t} \in T_{i}, d^{f} \in F_{i}} \mathbf{1}(\ln(d^{t}|c_{i}) > \ln(d^{f}|c_{i}))}{\sum_{i=0}^{n} \sum_{d^{t} \in T_{i}, d^{f} \in F_{i}} \mathbf{1}(\ln(d^{t}|c_{i}) < \ln(d^{f}|c_{i}))}$$
(2)

3.5 Experimental Setup

Our experiments require knowledge about the training data of LLMs. Therefore, we mainly study OLMo models (Team OLMo et al., 2024; Groeneveld et al., 2024) that provide publicly available pre- and post-training data. Additionally, we include Llama-3.1 (Meta AI, 2024) (Llama-3.1-Tulu-3-8B), which was post-trained with publicly available data. More model details are provided in Appendix A.3. We test on Natural Question (NQ) (Kwiatkowski et al., 2019), SicQ (Johannes Welbl, 2017), TrivialOA (Joshi et al., 2017), PopOA (Mallen et al., 2023) and TruthfulQA. Each of the first four datasets contains 1,000 randomly selected samples, and TruthfulQA (Lin et al., 2022) has 817. 174

175

176

177

178

179

180

181

183

184

185

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

209

210

211

4 **Result and Analysis**

Table 1 and Fig. 4 show the truthfulness results of different OLMo models on five benchmark datasets. We observe (1) a bigger model is not more truthful; (2) truthfulness is not correlated to task accuracy; (3) post-training schemes have limited impacts; pretraining schemes are more likely to determine the truthfulness of the model. Below is a more detailed analysis. We also provide ablation and case studies in Appendix A.6 and A.7 for more validations.

Which data samples are more influential for confidence generation? Confidence generation in OLMo2-13B Instruct (INS) is more influenced by confidence-related training data, while contentrelated data plays a greater role in OLMo-7B and OLMo2-7B INS models. As shown in Fig. 3, OLMo2-7B INS assigns high influence scores to samples that are semantically aligned with the question content. In contrast, OLMo2-13B INS is impacted by unrelated samples that contain superficial cues, such as the keyword "confidence".

Do LLM truthfulness relate to task accuracy? We hypothesize that high task accuracy in LLMs

Search	Instruct (INS)		NQ	S	SciQ	Tri	viaQA	Trut	hfulQA	Po	opQA	A	A11
Data	Model	ccr	acc	ccr	acc	ccr	acc	ccr	acc	ccr	acc	ccr	acc
Pre Post Pre+Post	OLMo2-13B	0.76 0.78 0.77	29.50	0.71 0.85 0.78	57.50	0.84 0.85 0.84	39.30	0.73 0.80 0.77	16.40	0.74 0.74 0.74	21.00	0.75 0.80 0.78	33.36
Pre Post Pre+Post	OLMo2-7B	1.08 1.85 1.43	23.60	1.01 1.63 1.32	50.50	1.17 1.73 1.45	32.20	1.43 1.67 1.56	15.42	1.02 1.74 1.37	16.10	1.12 1.72 1.42	28.03
Pre Post Pre+Post	OLMo-7B	1.22 1.28 1.25	17.30	0.95 1.17 1.05	39.80	1.08 1.29 1.18	27.00	1.30 1.25 1.27	20.93	1.06 1.15 1.10	11.90	1.11 1.22 1.16	23.48
	Average	1.14	23.47	1.07	49.27	1.23	32.83	1.20	17.58	1.07	16.33	1.14	28.62

Table 1: Result of LLMs' truthfulness measured with **ccr** and accuracy (**acc**) on different datasets. The training data used for estimating the influence score is searched from pre-training (Pre) data, post-training (Post) data or a combination of both (Pre+Post). We cross out the results that are not significant (p>0.05) in the mean influence score difference between the two comparison sets. Findings are robust across different training data settings.



Figure 4: Impact of different post-training schemes on the truthfulness assessment. We plot the training stages of the same model over time. The examined training samples are from the post-training data of the corresponding model. We include a new model, Llama3-8B, which has been post-trained with similar data to OLMo models.

may reflect exposure to similar examples during training, leading the model to rely on familiar, content-relevant information when generating confidence. If this is the case, we expect a positive correlation between task accuracy and truthfulness. However, as shown in Table 1, this correlation does not hold: the SciQ dataset achieves the highest average accuracy but a relatively low truthfulness score, while datasets like TruthfulQA, despite lower accuracies, exhibit higher truthfulness, suggesting LLMs may not rely on relevant training data that have seen during training to estimate confidence.

Do post-training schemes impact LLM truthfulness? The results in Fig. 4 demonstrate that post-training schemes can have opposite effects on different models. For example, DPO or RLVR improves the truthfulness for OLMo-7B and Llama3-8B but not for OLMo2-13B and OLMo2-7B models. We also observe that post-training schemes generally do not reverse the truthfulness results, i.e., scores are either all under one or above one, which indicates that pre-training schemes are more likely to determine the truthfulness.

5 Conclusion

LLMs often exhibit overconfidence, raising the question whether their confidence is grounded in truthful, content-relevant training data. This paper introduces TracVC, a method for tracing the origins of LLM confidence back to specific types of training samples. Our analysis reveals that OLMo2-13B models are more influenced by confidence-related, rather than content-related data when estimating confidence, suggesting a risk of trustworthiness in their expressed certainty. We further find that truthfulness does not correlate with model size, task accuracy or the extent of post-training. These findings motivate future work on how pre-training schemes shape confidence grounding, with the goal of building models whose confidence more reliably reflects truthful reasoning. Finally, our proposed TracVC method can be extended to interpret different types of model outputs from a data-driven perspective.

235

236

237

238

240

241

242

243

245

246

247

248

249

251

253

254

224

233

212

213

214

215

255

259

261

262

264

265

266

270

271

272

274

277

278

280

282

284

286

290

291

293

294

295

296

297

298

301

302

304

305

306

Limitations

Our findings are constrained to a limited set of LLMs, primarily the OLMo and Llama families, due to restricted access to the pre-training and posttraining data of other proprietary models. Since TracVC relies on analyzing training data influence, the lack of transparency and availability of training data for widely used commercial LLMs (e.g., GPT-4, Claude, Gemini) limits the generalizability of our conclusions. Future work could extend this analysis as more open-source models and datasets become available.

Our proposed method, TracVC, is broadly applicable and can be extended to trace various types of model outputs beyond confidence, such as specific answer choices. However, in this work, we focus solely on validating the method in the context of verbalized confidence. This focused scope allows for a clearer analysis of our method's effectiveness. Future work can extend our method to explore how the training data influences other forms of LLM output.

References

- Irina Bejan, Artem Sokolov, and Katja Filippova. 2023. Make every example count: On the stability and utility of self-influence for learning from noisy NLP datasets. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10107–10121, Singapore. Association for Computational Linguistics.
 - Tyler A. Chang, Dheeraj Rajagopal, Tolga Bolukbasi, Lucas Dixon, and Ian Tenney. 2024. Scalable Influence and Fact Tracing for Large Language Model Pretraining. *arXiv preprint*. ArXiv:2410.17413 [cs].
 - Jiuhai Chen and Jonas Mueller. 2024. Quantifying uncertainty in answers from any language model and enhancing their trustworthiness. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5186–5200, Bangkok, Thailand. Association for Computational Linguistics.
 - Sang Keun Choe, Hwijeen Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung, Adithya Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, Jeff Schneider, Eduard Hovy, Roger Grosse, and Eric Xing. 2024. What is your data worth to gpt? Ilm-scale data valuation with influence functions. *Preprint*, arXiv:2405.13954.
- Elastic. 2025. Elasticsearch: Distributed, restful search and analytics engine. https://www.elastic.co/ elasticsearch/. Version 8.17.2, released February 11, 2025.

Yanai Elazar, Akshita Bhagia, Ian Magnusson, Abhilasha Ravichander, Dustin Schwenk, Alane Suhr, Evan Pete Walsh, Dirk Groeneveld, Luca Soldaini, Sameer Singh, Hannaneh Hajishirzi, Noah A. Smith, and Jesse Dodge. 2024. What's in my big data? In *ICLR*.

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

325

326

327

328

329

330

331

332

334

335

337

338

339

340

341

342

344

345

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Olmo: Accelerating the science of language models. *Preprint*.
- Zayd Hammoudeh and Daniel Lowd. 2022. Identifying a Training-Set Attack's Target Using Renormalized Influence Estimation. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, CCS '22, pages 1367–1381, New York, NY, USA. Association for Computing Machinery.
- Zayd Hammoudeh and Daniel Lowd. 2024. Training data influence analysis and estimation: a survey. *Machine Learning*, 113(5):2351–2403.
- Matt Gardner Johannes Welbl, Nelson F. Liu. 2017. Crowdsourcing multiple choice science questions. *arXiv:1707.06209v1*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611.
- Chethan Kadavath, Amanda Askell, Jackson Kernion, Tom Henighan, Ben Mann, Gretchen Krueger, Sarah Kreps, Aaron McKane, Gaurav Mistry, Joe Kim, et al. 2023. Prompting GPT-3 to be reliable. In Proceedings of the 11th International Conference on Learning Representations (ICLR).
- Pang Wei Koh and Percy Liang. 2017. Understanding Black-box Predictions via Influence Functions. In Proceedings of the 34th International Conference on Machine Learning, pages 1885–1894. PMLR. ISSN: 2640-3498.
- Abhishek Kumar, Robert Morabito, Sanzhar Umbet, Jad Kabbara, and Ali Emami. 2024. Confidence under the hood: An investigation into the confidenceprobability alignment in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1:*

365 366 367 Long Papers), pages 315–334, Bangkok, Thailand.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-

field, Michael Collins, Ankur Parikh, Chris Alberti,

Danielle Epstein, Illia Polosukhin, Jacob Devlin, Ken-

ton Lee, et al. 2019. Natural questions: a benchmark

for question answering research. Transactions of the

Association for Computational Linguistics, 7:453-

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.

TruthfulQA: Measuring how models mimic human

falsehoods. In Proceedings of the 60th Annual Meet-

ing of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3214–3252, Dublin,

Ireland. Association for Computational Linguistics.

Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das,

Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating

effectiveness of parametric and non-parametric mem-

ories. In Proceedings of the 61st Annual Meeting of

the Association for Computational Linguistics (Vol-

ume 1: Long Papers), pages 9802-9822, Toronto,

Canada. Association for Computational Linguistics.

capable openly available llm to date. Accessed: 2025-

Meta AI. 2024. Introducing meta llama 3: The most

Shiyu Ni, Keping Bi, Lulu Yu, and Jiafeng Guo. 2025.

Are large language models more honest in their probabilistic or verbalized confidence? In *Information*

Retrieval, pages 124–135, Singapore. Springer Na-

Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guil-

laume Leclerc, and Aleksander Madry. 2023. TRAK:

Attributing Model Behavior at Scale. In Proceedings

of the 40th International Conference on Machine

Learning, pages 27074-27113. PMLR. ISSN: 2640-

Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund

Sundararajan. 2020. Estimating Training Data Influence by Tracing Gradient Descent. In Advances in

Neural Information Processing Systems, volume 33,

Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groen-

eveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yul-

ing Gu, Shengyi Huang, Matt Jordan, Nathan Lam-

bert, Dustin Schwenk, Oyvind Tafjord, Taira An-

derson, David Atkinson, Faeze Brahman, Christo-

pher Clark, Pradeep Dasigi, Nouha Dziri, Michal

Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng

Liu, Saumya Malik, William Merrill, Lester James V.

Miranda, Jacob Morrison, Tyler Murray, Crystal

Nam, Valentina Pyatkin, Aman Rangapur, Michael

Schmitz, Sam Skjonsberg, David Wadden, Christo-

pher Wilhelm, Michael Wilson, Luke Zettlemoyer,

Ali Farhadi, Noah A. Smith, and Hannaneh Ha-

jishirzi. 2024. 2 olmo 2 furious. arXiv preprint

arXiv:2501.00656.

pages 19920-19930. Curran Associates, Inc.

Association for Computational Linguistics.

- 36 36
- 37 37

466.

05-01.

3498.

ture Singapore.

- 373 374 375
- 3

379

- 333
- 384 385 386
- 33
- 390
- 39 39
- 395 396
- 390 399 400
- 401 402 403 404 405
- 406 407
- 408 409 410
- 411 412
- 412 413 414

415 416

- 416 417 418 419
- 420

Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics. 421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

- vLLM Contributors. 2023. vllm: A high-throughput and memory-efficient llm serving library. https: //github.com/vllm-project/vllm. Accessed: 2025-01.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. LESS: Selecting Influential Data for Targeted Instruction Tuning. *arXiv preprint*. ArXiv:2402.04333 [cs].
- Yuxi Xia, Pedro Henrique Luz de Araujo, Klim Zaporojets, and Benjamin Roth. 2025. Influences on llm calibration: A study of response agreement, loss functions, and prompt styles. *Preprint*, arXiv:2501.03991.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. In *Proceedings* of the 12th International Conference on Learning Representations (ICLR).
- Chih-Kuan Yeh, Ankur Taly, Mukund Sundararajan, Frederick Liu, and Pradeep Ravikumar. 2022. First is Better Than Last for Language Data Influence. *arXiv preprint*. ArXiv:2202.11844 [cs].
- Kaitlyn Zhou, Jena Hwang, Xiang Ren, and Maarten Sap. 2024. Relying on the unreliable: The impact of language models' reluctance to express uncertainty. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3623–3643, Bangkok, Thailand. Association for Computational Linguistics.

A Appendix

A.1 Prompt Details

First-stage prompt

This prompt aims to get the LLM answer of a question.

User: Answer the question, give ONLY the answer, no other words or explanation: Where does most of our food come from? $\rightarrow q$

Expected output from LLM: agriculture

Second-stage prompt

We use the LLM answer from the first stage prompt to generate this prompt in order to get the model's confidence on its answer. As we want to investigate the influences of confidence, therefore,

566

567

we ensure this stage of prompt targets to guide the LLM to generate only confidence.

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

502

504 505

507

User: Answer the question, give ONLY the answer, no other words or explanation: Where does most of our food come from?

Assistant: **agriculture** $\rightarrow a$ (target answer is angiosperms, but we do not care about the correctness of the answer in our setting.)

User: Provide the probability that your answer is correct. Give ONLY the probability between 0.0 and 1.0, no other words or explanation. $\rightarrow p$

Expected output from LLM: 0.9

A.2 Training Data Influence Estimation

The method for attributing training data influence in this work is inspired by TracIn (Pruthi et al., 2020), which, in its original formulation, aims to estimate the influence that training on an instance z had, when predicting a test instance z': It does so by measuring the similarity between the gradients of the model's loss function, when evaluated on zand z' respectively, w.r.t some set of parameters w_t , at a series of checkpoints T:

$$\phi_{\text{TracInCP}}(z, z') = \sum_{\forall t \in T} \eta_t \nabla \ell(w_t, z) \cdot \nabla \ell(w_t, z')$$
(3)

However, in line with previous work, we utilize cosine similarity for our content-over-confidence ratio, to reduce the impact of gradient magnitude on our comparison (Hammoudeh and Lowd, 2022, 2024; Park et al., 2023; Xia et al., 2024).

We compute the gradient with respect to the model's input embeddings. Note that this score captures information about model behavior beyond the word embedding layer, as the gradient chains through higher layers as well (Yeh et al., 2022).

A.3 Detailed Experimental Setup

Our studied OLMo models span different model 509 sizes and versions: OLMo2-13B (OLMo-2-1124-510 7B-Instruct), OLMo2-7B (OLMo-2-1124-13B-Instruct), OLMo-7B (OLMo-7B-Instruct-hf). Ad-512 ditionally, we study how different post-training 513 schemes impact the ccr score. Thus, we in-514 clude Llama-3.1 (Meta AI, 2024) (Llama-3.1-Tulu-515 516 3-8B) that has the post-training data accessible. For each model, we compute the ccr scores for 517 checkpoints after post-training with supervised-518 fine-tuning (SFT), direct preference optimization 519 (DPO), and reinforcement learning with verified 520

reward (RLVR) (the last training step for the instruction model). We employ the vLLM (vLLM Contributors, 2023) library for LLM inference and serving. All LLMs are set with temperature equal to 0 to ensure consistent outputs. All our experiments are conducted on NVIDIA HGX H100, requiring a total of approximately 642 GPU hours for testing various scenarios.

A.4 Accuracy of More Models

In addition to the accuracy performance of models shown in Table 1. We report the accuracy of models from different post-training stages in Fig. 5.

A.5 Details of the Search Data

We show the details of all the pre-training and posttraining data in Table 2. We use dolma v1.7 as pre-training data for OLMo models, which is provided as a publicly available online search engine by Elazar et al. (2024).

A.6 Ablation Study

We validate our method in three additional settings:

Influence Estimation with Verbalized Confidence. In the setup described in the paper, we estimate the influence of a training sample d on the model's generated confidence c by computing the cosine similarity between two gradients of the model's loss function: one taken when evaluated on the training sample d, and the other when evaluated on (q, a, p). For completeness, we also provide results for (q, a, p, c), i.e., with the verbalized confidence c included in the input in Table 3. Our influence score is then defined as:

$$\operatorname{Inf}(d|c_i) = \frac{\nabla \ell(w_i, d) \cdot \nabla \ell(w_i, (q_i, a_i, p, c_i))}{\|\nabla \ell(w_i, d)\| \cdot \|\nabla \ell(w_i, (q_i, a_i, p, c_i))\|}$$
(4)

The results of the setting are shown in Table 3. We observe similar scores to Table 1 in the main paper. However, we find that more result scores are not significant, i.e., Table 1 only has one insignificant score. This indicates that **our method presented in the main paper is more robust in providing significant results**.

Random Baseline. To further validate our method, we apply TracVC to a random baseline. Specifically, we include a setting where we replace the confidence-related set with random samples from different datasets, resulting in a **r**andom set $R = \{d_1^r, ..., d_{10}^r\}$. A corresponding metric similar to **ccr** is named as *Content-over-Random* ratio

Model	Pre-training Data	Post-training Data
OLMo2-13B INS	dolma v1.7	tulu-3-sft-olmo-2-mixture, olmo-2-1124-13b-preference-mix, RLVR-MATH
OLMo2-13B-DPO	dolma v1.7	tulu-3-sft-olmo-2-mixture, olmo-2-1124-13b-preference-mix
OLMo2-13B-SFT	dolma v1.7	tulu-3-sft-olmo-2-mixture
Llama3-8B-INS	-	tulu-3-sft-mixture, RLVR-GSM-MATH-IF-Mixed-Constraints
Llama3-8B-DPO	-	tulu-3-sft-mixture, llama-3.1-tulu-3-8b-preference-mixture, llama-3.1-tulu-3-8b-
		preference-mixture
Llama3-8B-SFT	-	tulu-3-sft-mixture
OLMo2-7B-INS	dolma v1.7	tulu-3-sft-olmo-2-mixture, olmo-2-1124-7b-preference-mix, RLVR-GSM
OLMo2-7B-DPO	dolma v1.7	tulu-3-sft-olmo-2-mixture, olmo-2-1124-7b-preference-mix
OLMo2-7B-SFT	dolma v1.7	tulu-3-sft-olmo-2-mixture
OLMo-7B-INS	dolma v1.7	tulu-v2-sft-mixture, ultrafeedback-binarized-cleaned
OLMo-7B-SFT	dolma v1.7	tulu-v2-sft-mixture

Table 2: Details of the search data. We use dolma v1.7 as a subset of the full pre-training data of OLMo-2-13B and OLMo-2-7B models.

Search Data	Model	NQ	SciQ	TriviaQA	TruthfulQA	PopQA	All
Pre	OLMo2-13B-INS	0.75	0.71	0.82	0.70	0.73	0.74
Post	OLMo2-13B-INS	0.73	0.85	0.79	0.70	0.71	0.75
Pre+Post	OLMo2-13B-INS	0.75	0.78	0.81	0.70	0.72	0.75
Pre	OLMo2-7B-INS	1.07	1.02	1.22	1.54	1.06	1.15
Post	OLMo2-7B-INS	1.95	1.52	1.87	1.84	1.71	1.77
Pre+Post	OLMo2-7B-INS	1.43	1.25	1.49	1.65	1.32	1.41
Pre	OLMo-7B-INS	1.22	.0.96	1.08	1.29	1.06	1.11
Post	OLMo-7B-INS	1.30	1.18	1.34	1.33	1.18	1.26
Pre+Post	OLMo-7B-INS	1.25	1.05	1.19	1.30	1.11	1.17

Table 3: Result of LLMs' truthfulness measured with **ccr** on different datasets **when including the verbalized confidence in the test samples** for influence estimation. The training data used for estimating the influence score is searched from pre-training (Pre) data, post-training (Post) data or a combination of both (Pre+Post). We cross the results that are not significant (p>0.05) in the mean influence score difference between the two comparison sets. Findings are robust across different training data settings and similar to the ones shown in Table 1, but with more insignificant values on SciQ dataset.

(crr):

568

569

570

571

572

573

574

575

576

577

578

582

583

584

585

$$crr = \frac{\sum_{i=0}^{n} \sum_{d^{t} \in T_{i}, d^{r} \in R_{i}} \mathbf{1}(\ln(d^{t}|c_{i}) > \ln(d^{r}|c_{i}))}{\sum_{i=0}^{n} \sum_{d^{t} \in T_{i}, d^{r} \in R_{i}} \mathbf{1}(\ln(d^{t}|c_{i}) < \ln(d^{r}|c_{i}))}$$
(5)

A higher **crr** score indicates that LLMs are more influenced by content-related training data than random training samples. The evaluated results shown in Table 4 suggest that **LLMs are generally more impacted by content-related training data than random variants.**

Identical Training Data to Input. In this setting, we assume there is an ideal, identical training sample, the same as the input query. We study whether such identical training samples could provide a higher truthfulness score for the tested models. The formula of the Identical-Content over Identical-Confidence ratio (**icicr**) is defined as:

$$icicr = \frac{\sum_{i=0}^{n} \mathbf{1}(\ln(q_i, a_i | c_i) > \ln(p, c_i | c_i))}{\sum_{i=0}^{n} \mathbf{1}(\ln(q_i, a_i | c_i) < \ln(p, c_i | c_i))}$$
(6)

The results are shown in Table 5. We observe

many insignificant scores, and the scores across different post-training checkpoints of one model conflict with each other (OLMo2-7B models). These findings strongly indicate that **it is not effective to use identical training samples for investigating our research question** in the main paper.

A.7 Case Study

We provide more detailed case studies in Table 11-10 for the most influential data samples we find by searching in pre-training data of LLMs for each test dataset. Table 11 represents the most influential samples we retrieved from post-training data and their corresponding influence scores measured with TracVC.

599

586

Search Data	Model	NQ	SciQ	TriviaQA	TruthfulQA	PopQA	All
Pre	OLMo2-13B-INS	1.13	0.89	1.12	0.90	1.09	1.03
Post	OLMo2-13B-INS	1.63	1.71	1.62	1.76	1.72	1.68
Pre	OLMo2-7B-INS	0.97	1.19	0.90	1.12	1.22	1.07
Post	OLMo2-7B-INS	1.30	1.04	1.09	1.11	1.17	1.14
Pre	OLMo-7B-INS	1.34	0.78	0.87	1.26	.0.98	1.02
Post	OLMo-7B-INS	1.15	1.04	1.11	1.09	1.02	1.08

Table 4: Result of LLMs' truthfulness measured with **Content-over-Ranodom ratio** on different datasets. The training data used for estimating the influence score is searched from pre-training (Pre) data or post-training (Post) data. We cross the results that are not significant (p>0.05) in the mean influence score difference between the two comparison sets.

Search Data	Model	NQ	SciQ	TriviaQA	TruthfulQA	PopQA	All
Identical	OLMo2-13B-INS	0.19	0.20	0.27	0.19	0.14	0.20
Identical	OLMo2-13B-DPO	.0.28	.0.37	0.42	0.34	0.18	0.31
Identical	OLMo2-13B-SFT	0.11	0.12	0.12	0.14	.0.13	0.12
Identical	Llama3-8B-INS	0.98	1.00	0.92	1.28	1.65	1.13
Identical	Llama3-8B-DPO	1. 0 7	1.00	0.95	1.29	1.67	1.16
Identical	Llama3-8B-SFT	1.11	1.15	0.93	1.30	1.82	1.22
Identical	OLMo2-7B-INS	0.31	0.29	0.23	0.21	0.12	0.23
Identical	OLMo2-7B-DPO	0.94	0.83	0.96	0.75	1.05	0.91
Identical	OLMo2-7B-SFT	4.00	2.05	2.82	1.08	1.74	2.11
Identical	OLMo-7B-INS	1.42	1.38	1.77	1.49	1.92	1.58
Identical	OLMo-7B-SFT	4.05	5.80	5.45	4.63	9.53	5.51

Table 5: Result of truthfulness measured with **Identical-Content over Identical-Confidence ratio**. The examined data for the influence score is identical to the input query. We cross out the results that are not significant (p>0.05) in the mean influence score difference between the two comparison sets.



Figure 5: Accuracy of different posting-training checkpoints. Different color of lines presents different datasets like Fig. 4 in the main paper.

Model	Answer	Confidence	Max Inf. Score	Most Influencial Document from Pre-training Data
OLMo2- 13B-INS	John An- glin	0.8	0.01	[Confidence-related]: Type of the result message. Transcript text representing the words that the user spoke. Populated if and only if messagetype equals TRANSCRIPT. If false, the StreamingRecognitionResult represents an interim result that may change. If true, the recognizer will not return any further hypotheses about this piece of the audio. May only be populated for messagetype = TRANSCRIPT. The Speech confidence between 0.0 and 1.0 for the current portion of audio. A higher number indicates an estimated greater likelihood that the recognized words are correct. The default of 0.0 is a sentinel value indicating that confidence was not set. This field is typically only provided if isfinal is true and you should not rely on it being accurate or even set. An estimate of the likelihood that the speech recognizer will not change its guess about this interim recognition result: - If the value is unspecified or 0.0, Dialogflow didn't compute the stability. In particular, Dialogflow will only provide stability for TRANSCRIPT results with isfinal = false Otherwise, the value is in (0.0, 1.0] where 0.0 means completely unstable and 1.0 means completely stable. Word-specific information for the words recognized by Speech in transcript. Populated if and only if messagetype = TRANSCRIPT and [InputAudioConfig.enablewordinfo] is set. Time offset of the end of this Speech recognition result relative to the beginning of the audio. Only populated for messagetype = TRANSCRIPT. DTMF digits. Populated if and only if
OLMo2- 7B-INS	Frank Lee Morris	0.9	0.20	IContent–related]: Only the Morris family knows the real story of the escape from Alcatraz Island by Frank Morris and the Anglin Brothers. But now the story will be told. I believe that the man in this video Bud Morris could be the real Frank Morris. What do you all think??? Did Frank Morris Survive After He Escaped From Alcatraz? Alcatraz PrisonThe cell of Frank Lee Morris. Is this the picture that 'proves' John and Clarence Anglin DID escape Alcatraz? CBS station KPIX-TV obtained a letter allegedly written by one of the men who escaped from the federal prison on San Francisco's Alcatraz Island in 1962. VIEWER DISCRETION IS ADVISED. With mastermind Frank Morris at the helm, three impacts at Alextrag arean suppagefully account from their incorporation.
OLMo- 7B-INS	The per- son who escaped from Alcatraz, on June 11, 1962, was Frank Morris & the Three Stooges; though none of their bodies were ever found, it is gen- erally accepted that they were re- sponsible for mak- ing their escape.	1.0 - The proba- bility I accu- rately remem- bered the details of the Alcatraz escape and iden- tified the correct perpetra- tors.	0.22	Icontent-related): 17 items matched your search for "Frank Lee Morris" Frank Morris was smart enough to physically ready himself for the one-mile swim over the six months or more he spent preparing for his bid for freedom Frank Morris was smart (with a tested IQ of 133), too Who are California's most wanted criminals? Crime Scene Kurtis Alexander Men like Frank Lee Morris, who escaped from Alcatraz in 1962, and suspected eco-terrorist and Berkeley native Daniel Andreas San Diego have long eluded the grasp of the law The Cartier Hunter: Wanted by the US Marshals Service, Oakland Police Department, and the California Department of Corrections Fugitive Apprehension Team for murder and parole violations Frank Lee Morris, one of the Alcatraz escapees of 1962. Email EMAILADDRESS Twitter: EvanSernoffsky Joe Tate works on a jail doors of former Alcatraz prisoner Frank Lee Morris, who escaped from Alcatraz in 1962. If there was ever an inmate who was 55 years ago, Alcatraz guards realized three inmates escaped their cells. They've never been found. Their morning inmate count revealed there missing men: Frank Morris and brothers Clarence and John Anglin And five, including Frank Morris, John Anglin and Clarence Anglin, were never found Frank Lee The two, Dari Lee Parker and John Paul Scott, were missed at 5:47 p.m Escapers such as Frank Lee Morris and John and Clarence Anglin, who are known to have entered the water last June. The three were identified as Frank Lee Morris, 35, from Louisiana; John Anglin, 32, and Clarence Anglin, 31, two of three Florida brothers doing time for an Alabama bank job. Their breakout was an incredible Joe Tate works on a jail doors of once famed Alcatraz poisoner Frank Lee Morris (dumby head in bed) who escaped from Alcatraz 1960. If there was ever an inmate who was destined to escape from Alcatraz. Three: You escaped the Rock in '62, and yesterday I said all was forgiven and you three (Frank Lee Morris who made a daring escape in 1962 and U

Table 6: The most influential document and its corresponding influence score for the confidence generation when testing an LLM on a question from NQ. The test question is "Who was the person who escaped from Alcatraz?"

Model	Anwer	Confidence	Max Infl. Score	Most Influencial Document from Pre-training Data
OLMo2- 13B-INS	nurture	0.9	0.01	[Content-related]: You are given a question, its answer, and a sentence that supports the question, i.e., the answer to the question is inferable from the sentence. In this task, you need to paraphrase the given sentence so that the paraphrased sentence still supports the question i.e. you can still infer the answer to the question from the paraphrased sentence. Do not write a paraphrase with a minor change in the given sentence e.g. replacing the word "one" with "a". Instead, try to write a paraphrase that contains new words, i.e. the words that are not present in the input sentence. Q: Question: Animal behavior can be said to be controlled by genetics and experiences, also known as nature and nurture. A: Animal behavior is controlled by both genetics and experience, otherwise referred to as nature and nurture.
OLMo2- 7B-INS	nurture	0.9	0.16	[Confidence-related]: <issue start=""><issue comment="">Title: blob.sentiment doesn't return polarity for "Erfolg" (or any other noun I tried) username 0: Hi there, thanks a lot for this great package! While analyzing some texts with TextBlob-DE I stumbled upon the following difference between TextBlob and TextBlob-DE as TextBlob from textblob import TextBlob as TextB blob = TextBlob/DE as TextBlob from textblob import TextBlob as TextB blob = TextBlob/DE as textBlob.pl = TextBlob/DE set textBlob-DE seems to be unable to find the polarity for "Erfolg" in the German sentiment file: <word confidence="1.0" form="Erfolg" intensity="1.0" polarity="1.0" pos="NN" subjectivity="0.0"></word> On the other hand TextBlob is obviously able to include the polarity for "success" into the sentiment calculation: <word confidence="0.9" form="success" id="n-14474894" intensity="1.0" polarity='0.3"' pos="NN" sense=" a state of prosperity or fame" subjectivity="0.0" wordnet=""></word> I would very much appreciate if you could give me any advice on how to resolve this issue. Kind regards, Andreas cissue comment_username 0: I fixed this problem in my local copy by changing the following statement in Sentiment.load(): 'w.attrib.get("form"),' changed to 'w.attrib.get("form").lower(),' Now the lowercased words in the input text can be matched with the lowercased words in the loaded lexicon. Here the results of a "tiny test": Text; Kommentar; sentiment(polarity) Das ist ein Erfolg; positives Substantiv; 1.0 Das war ein kein Erfolg; negiertes positives Substantiv; -0.5 Der Text with FoxtBlob-DE I stumbled upon the following difference between TextBlob and TextBlob-DE concerning the polarity of nouns as returned by blob.sentiment) "' This code print (blob.sentiment) blob = TextBlob'.DE is a success') print(blob.sentiment) in the sentiment file: <word "tiny="" ,="" a="" be="" can="" confidence="1.0" form="Erfolg" here="" in="" input="" intensity="0.0" kommentar;="" lexicon.="" loaded="" lower-cased="" lowercased="" matched="" now="" of="" results="" sentim<="" td="" test":="" text="" text;="" the="" with="" words=""></word></issue></issue>
OLMo- 7B-INS	Nurture.	I cannot directly provide a proba- bility, but the state- ment is accurate. Accord- ing to scientific research, genetics (nature) and expe- riences (nurture) both influence animal behavior, with genetics typically having a more sig- nificant and stable impact (higher probabil- ity)	0.13	Iconfidence-related]: I am using Random Forests, XGBoost and SVMs to classify whether the home team wins or the away team wins their bowl game (in college football). I trained the models on all the games during the season. I've come across something that is a bit weird and can't explain. I calculated a prediction confidence by subtracting the class probabilities. The XGBoost confidence values are consistency higher than both Random Forests and SVM's. I've attached the image below. I did some hyper-parameter tuning for all of my models and used the best parameters based on testing accuracy. minimum split criteria of 5 rows. I wasn't clear with my question: Why exactly does XGBoost prefer one class greatly to the other? In comparison to these other methods. I'm trying to figure out why my prediction confidences of a class are so high for XGboost . Rather than answering why XGBoost give very confident predictions, I will answer why random forest and SVM give not-so-confident predictions. Random forest probability estimates are given by the percentage of the forest that predicted a particular class. For example, if you have 100 trees in your forest and 81 of them predict some class for some example, the probability estimate for that example belonging to that class is calculated to be $\frac{81}{100} = 0.81$. Because of the random nature of the ensemble members, it's very unlikely that each individual tree will end up with the correct prediction, even if the majority do. This makes probability estimates from random forests shy away from the extreme ends of the scale. SVM is a slightly different case, because they are unable to produce probability estimates directly. Typically, Platt scaling (essentially logistic regression) is used to scale the SVM output to a probability estimate. This has the added benefit of calibrating the probability estimates, meaning the prediction, it actually has approximately a 80% chance of being correct. For a problem like this where there's a lot of noise (underdog teams do win sometime

Table 7: The most influential document and its corresponding influence score for the confidence generation when testing an LLM on a question from Sciq. The test question is "Animal behavior can be said to be controlled by genetics and experiences, also known as nature and what?"

Model	Anwer	Confidence	Max Infl. Score	Most Influencial Document from Pre-training Data
OLMo2- 13B-INS	penultimate	0.99	0.06	[Content-related]: "Penultimate" comes from a Latin word that means "almost ultimate," so the next to last book in a series, the next to last day of a vacation, and the next to last game in a player's career are all penultimate items or events. "Penultimate" is not the best of the bunch or the last of something; it is the second-best of the bunch or second-to-the-last of something. Believe me, ladies and gentlemen, there is nothing penultimate about this one. This, ladies and gentlemen, is the proverbial it. After this, there is void emptiness oblivion absolute nothing. "Penultimate" was actually a noun before it became an adjective. According to the Online Etymology Dictionary, "penultimate" referred to the "next to the last syllable of a word or verse." For example, I found an old dictionary from the 1800s that instructed people to "accent the penultimate" when explaining how to pronounce Greek and Latin proper names. The Latin prefix "panee." (shortened here to "pen-") means "almost" or "nearly." It's not very common anymore. Most words that use it now are obscure or rare (for example, "peneseismic" means regions where earthquakes occur only rarely or only of small magnitude, so it means something like "nearly seismic"), but one word still in common use is "peninsula," which means "almost island." Another word from the same root that you might have heard, especially if you have watched an eclipse, is "penumbra." "Umbra" means "shade or shadow," so a penumbra is almost a shadow or a partly shaded area. During a total solar eclipse, the total eclipse is only visible from certain parts of earth that are properly aligned to see it. Those people are covered by the "umbra"—the shadow. But people outside that region still see a partial eclipse, and they are said to be covered by the "penumbra"—the partial shadow. So the next time you want to describe something that is the best, simply call it the best or the
OLMo2- 7B-INS	penultimate	0.99	0.34	ultimate—the ultimate prize in the ratile—not the penultimate prize. [Confidence-related]: Q: Customising matplotlib cmaps I have some normalised histogram data in array of shape (12,1): \gg hnorm array([[0,], [0,], [0, 01183432], [0.0295858], [0.04142012], [0.04142012], [0.04142012], [0.03550296], [0.01775148], [1,], [0.98816568], [0.56213018], [0,]]) I'd like to plot this in 'heatmap' style. I am doing this like so: import matplotlib.pyplot as plt plt.imshow(hnorm, cmap='RdBu',origin='lower') This works (axis formatting aside). However, I'd like to customise the colormap to fade from white to Red. I have attempted: import matplotlib.colors as col cdict = 'red': ((0.0, 0.0, 0.0), (0.0, 0.0, 0.0), (0.0, 0.0, 0.0), (0.0, 0.0, 0.0), (0.0, 0.0, 0.0), (0.0, 0.0), (0.0, 0.0), (0.0, 0.0), (0.0, 0.0), (0.0, 0.0), (0.0, 0.0), (0.0, 0.0), "bue': ((0.0, 1.0, 1.0), (1.0, 0.0), (0.0, 0.0), (0.0, 0.0), 'bue': ((0.0, 1.0, 1.0), (1.0, 0.0)) my cmap = mpl.colors.LinearSegmentedColormap('my colormap', cdict) plt.imshow(np.random.rand(20,20), cmap=my cmap, origin='lower', interpolation='none') plt.colors/shrink=.75) Another example showing how the two color items a
OLMo- 7B-INS	If you're looking for a word that means second- to-last in a list or in terms of com- pletion, the word is "subse- quent." So, the answer to the question is "subse- quent."	1.0 - The proba- bility I would use to deter- mine the accuracy of the in- formation provided is not ap- plicable, as I am simply providing an answer based on my train- ing and knowl- edge. The probabil- ity of my answer being correct is 100%.	0.20	[Confidence-related]: <issue start=""><issue comment="">Title: How to compute labelling probability after prediction? username 0: Hi, Thanks a lot for this wonderful piece of work. I am trying to calculate CTC loss to compute labeling probability after prediction. Please guide if it is possible Thanks <issue comment="">username 1: Could you please elaborate on what you are try- ing to do? <issue comment="">username 0: Hi, I actually want to have a score which tells me with what confidence can i claim that my prediction for a handwritten text is correct?. <issue com-<br="">ment>username 1: Hi username 0, A handwriting recogniser can be evaluated with the Word Error Rate ([WER](https://en.wikipedia.org/wiki/Word error rate)) or Character Error Rate (CER). In our implemen- tation, we used SCLITE to calculate it. You can see it [here](https://github.com/awslabs/handwritten- text-recognition-for-apache-mxnet/blob/master/ocr/utils/sclite helper.py). I hope this answers your question. <issue comment="">username 0: Hi username 1, Thanks for the reply. I understand what you mean. But my problem is this. I do not have ground truth(actual text) with me. So whenever I predict a handwritten text,can i log a confidence score along with it, telling how confident I am about the prediction. Something like this- No 3 in https://towardsdatascience.com/faq-build-a-handwritten-text- recognition-system-using-tensorflow-27648fb18519 Thanks <issue comment="">username 2: Because of how the CTC collapsing work, any confidence score would be quite hard to compute and rely on, because of repetitions, multiple possible paths, etc. However i you want a score that might be helpful to you, you could softmax normalize the prediction of each time step, and use the average probability for the characters of the final path as an indication of confidence. For example, with [a,b,c,] as a dictionary, sequence length of 6. " a 0.1 0.0 1.0 0.0 to 0.0 0.0 2.0 1.0 0.0 c 0.7 0.8 0.0 0.1 0.0 0.0 0.0 0.0 1.0 0.0 1.0 0.0 9." " ccabb" => cab" after CTC collapsing (0.7+0.8+1.0+0.7+1</issue></issue></issue></issue></issue></issue></issue>

Table 8: The most influential document and its corresponding influence score for the confidence generation when testing an LLM on a question from TriviaQA. The test question is "**If ultimate means last, what word means second-to-last?**"

Model	Anwer	Confidence	Max Infl. Score	Most Influencial Document from Pre-training Data
OLMo2- 13B-INS	No harm	0.9	0.03	[Confidence-related]: <issue start=""><issue comment="">Title: Making Serializable reflective username 0: Now all 'Serializable' objects can optionally report their classes "reflectively". Prior to this PR, only the 'Container' class was reflective. This meant that when a 'Container' was embedded in another 'Serializable' object, we lost the ability to control whether or not the reflective attributes were included in the JSON when invoking 'serialize()' on the parent object. Now 'serialize()', 'write json()', and 'to str()' support a 'reflective=True/False' argument that controls this behavior recursively. I wanted this behavior because it's not appropriate to include "reflective" attributes when generating JSON that's shared with users. As an added benefit, one can now call 'Serializable.from json()' to load any arbitrary Serializable object that was serialized with 'reflective=True'. Example: "'py import eta.core.geometry as etag import eta.core.objects as etao tl = etag.RelativePoint(0, 0) br = etag.RelativePoint(1, 1) bb = etag.BoundingBox(tl, br) obj = etao.DetectedObject("car", bb, confidence=0.9, index=1) obj attr = etao.ObjectAttribute("make", "Honda", confidence=0.9) obj.add attribute(obj attr) # no reflection (de- fault) print obj # with reflection print obj.to str(reflective=True) "' No reflection: "ijson "label": "car", "bounding box": "bottom right": "y": 1.0, "x": 1.0, "top left": "y": 0.0, "x": 0.0, "confidence": 0.9, "index": 1, "attrs": ["category": "make", "label": "car", "bounding box": " CLS": "eta.core.geometry.RelativePoint[", "Y": 0.0, "x": 0.0, "confidence": 0.9, "index": 1, "attrs": " CLS": "eta.core.objects.ObjectAttributeContainer", "ATTR CLS": "eta.core.objects.ObjectAttribute", "attrs": [" CLS": "eta.core.objects.ObjectAttribute, "category": "make", "label": "Honda", "confidence": 0.9] "" 'make", "label": "Honda", "confidence": 0.9] "" 'make", "label": "Honda", "confidence": 0.9] "" <issue comment="">username 1: I am going to merge this now.</issue></issue></issue>
OLMo2- 7B-INS	digestive discom- fort	0.8	0.29	[Content-related]: What Really Happens If You Eat Watermelon Seeds? Watermelon is a summer essential. What happens if you eat the black or white watermelon seeds? Are you in danger? Short Answer: No, you'll survive. Long Explanation Answer: Fact- Swallowing a watermelon seed will not cause a watermelon to grow in your stomach. When you swallow watermelon seeds raw, they move through your digestive tract without being digested. That's it. According to Spoon University, you can actually prepare watermelon seeds correctly to eat because they are full of many health benefits. These seeds are packed with protein, vitamins, and minerals. About 1/8 of a cup of watermelon seeds, there is a little bit of labor required, as they need to be sprouted, shelled, and dried." So remember, watermelon seeds completely safe to swallow, and they actually have many health benefits. If you eat them raw, you will miss out on all of these benefits.
OLMo- 7B-INS	If you eat wa- termelon seeds, you may expe- rience discom- fort and in- testinal block- age. Water- melon seeds contain hard, indi- gestible shells that must be broken down before they can be passed from the body.	0.0 - Cannot be deter- mined with available infor- mation. Water- melon seeds do not naturally occur in human diets, so there is no estab- lished data on the likeli- hood of adverse effects.	0.17	[Content-related]: How Much Watermelon Rind to Give to Your Dog? Can Dogs Eat Watermelon With Seeds? How Many Watermelon Seeds to Give to Your Dog? Why Do Dogs Like Watermelon Dogs Eat Black Watermelon Seeds? But you might wonder if it is safe for the dog. Eating watermelon is very healthy as it keeps the body hydrated. But when consumed in large amounts, watermelon and cause diarrhea, both in dogs and humans. Dogs can eat watermelon but with a few precautions. One might ask can dogs eat watermelon is very healthy as it loogs and numans. Dogs can eat watermelon but with a few precautions. One might ask can dogs eat watermelon yet dog are not indigestible and it can cause gastrointestinal blockage. Also, the hard rind can also damage your dogs' gums and teeth. Watermelon is very healthy and nutritious for dogs, but some parts can cause severe damage. It will cause an intestinal blockage that is not only painful but can also lead to surgery if not taken care of properly. So we will discuss in detail whether or not your dog can eart rind or seed, or watermelon as it depends on various other facts addressed in detail in the following sections. The rind of the watermelon is not safe to eat. The rind has two portions, the light green inside part and the hard outside. Many dogs can pick over the light green portion, but that part is very firm and hard to chew. If any of such hard parts are taken in by the dog. then the dog many not chew it thoroughly and swallow it as a whole. This is a point that can cause issues. Any such indigestible, so that it can lead to constipation, vomiting, lethargy, and abdominal pain in the dog. The smaller dogs are at greater risk in this condition. They have a small stomach, and their teeth are not that strong enough to thoroughly the whet rind. They have a small stomach, and their teeth are not were than and unso the watermelon as in depared on a waterget or is to chew the hard rind. Larger dogs have the strength and larger body cavities, so they can chew the rind a bin and may be a

Table 9: The most influential document and its corresponding influence score for the confidence generation when testing an LLM on a question from TruthfulQA. The test question is "What happens to you if you eat watermelon seeds?". Document text abbreviated.

Model	Anwer	Confidence	Max Infl. Score	Most Influencial Document from Pre-training Data
OLMo2- 13B-INS	Mike Leigh	0.9	0.02	[Content-related]: This episode is Filmwax UK Edition! Mike Leigh gets the full retrospective treatment by Film at Lincoln Center with HUMAN CONDITIONS: THE FILMS OF MIKE LEIGH (5/27-6/8). He stops by to discuss. And the creative team behind a new comedy out of the UK, director Craig Roberts & screenwriter Simon Farnaby. The film is called "The Phantom of the Open" and it opens on June 3rd. Mike Leigh returns to Filmwax. He was last on Episode 547 back in 219 discussing his most recent film, "Peterloo". This time he appears as he is being given the full retrospective treatment at Film at Lincoln Center with Human Conditions: The Films of Mike Leigh. The retrospective, which takes place from May 27th through June 8th includes all of Leigh's features and a few other choice surprises. Leigh will be on hand for Q&A's this weekend. The director Craig Roberts and screenwriter Simon Farnaby have created a new comedy called "The Phantom of the Open" which opens Friday, June 3rd. The story follows Maurice Flitcroft (Mark Rylance), a dreamer and unrelenting optimist who managed to gain entry to The British Open Golf Championship Qualifying in 1976 and subsequently shot the worst round in Open history, becoming a folk hero in the process. The film also scalub Huwking & Phave Greace Tome is case
OLMo2- 7B-INS	Barry Hines	0.9	0.26	[Content-related]: Join The Showroom and And Other Stories for an evening celebrating the influential South Yorkshire-born novelist and screenwriter, Barry Hines, and the reissue of his masterpiece of nature writing and rural class conflict, The Gamekeeper. Born into a mining family in a village near Barnsley, Barry Hines (1939-2016) worked first in a coal mine before going to college, working as a teacher, and then becoming a full-time writer of fiction and screenplays for film and television. Hines is best known for A Kestrel for a Knave, a novel that has never been out of print in Britain and was filmed by Ken Loach as the widely acclaimed Kes. For over forty years he documented working-class lives with a boundless humanity, deep empathy, and ultimately, hope. Sheffield-based publisher AOS are proud to be reissuing a selection of Hines' novels over the coming years, including classics, lesser-known gems and until-now undiscovered work. This April, AOS published The Gamekeeper, Hines' gripping novel of rural working-class life through the changing seasons, seen through the eyes of a gamekeeper on a country estate in the North of England. To mark its publication, join us for a screening of the rarely-seen film adaptation of The Gamekeeper, adapted by Hines and directed by Ken Loach. There will also be a series of short talks on Barry Hines' work and legacy from those that he influenced and that knew and worked with him. There will be an introduction before this screening. David Forrest is Professor of Film and Television and Studies at the University of Sheffield. His most recent book is New Realism: Contemporary British Cinema (2020), and with Sue Vice he is the co-author of Barry Hines: Kes, Threads and Beyond (2017). Ron Roseis a playwright and scriptwriter, born in Sheffield and lives in Doncaster. He's had over 70 plays performed professionally at theatres across the country including two verbatim dramas set in the '84/'85 Miners' Strike Never the Same Again and The Enemy Within; the WW2 mining strikes dr
OLMo- 7B-INS	Richard Curtis was the screen- writer for Mean- time.	1.0 - I am able to pro- vide an accurate answer.	0.18	Rescue and Resistance (2021). [Confidence-related]: I am interested to find a way to identify if two sets of data can be considered statistically different at 95% Confidence level (or any other). To be more specific, my data sets are composed of 5 values. They correspond to 5 readings of the same detector by one microscope. So, two detectors with 5 readings each can represent the same value or not. The problem is to find a method to answer this question. I'm wondering if a t-test is the proper tool. Am I right? I would like to do the statistical test in R. I am slightly unsure what you are testing for. Are you testing for concordance (you expect the two sets of data to be the same) or are you testing to find a difference between the two sets? The confidence interval (-2.6, 1.4) includes zero, so the null hypothesis cannot be rejected and no difference is observed. Here, the p-value is 1.0, so a similar conclusion is made that the null cannot be rejected and no difference is observed. For Concordance However, if you are hypothesizing the more subtle outcome that the two sets of paired data concord with each other. The statistic of interest, rho, is similar to Pearson's product-moment correlation coefficient, but adjusted for exact agreement along the $x=y$ line (note, this is R, not R-square). Closer to 1.0 is strong concordance; closer to 0.0 is no concordance. In R, use the epiR package. Here, the rho is modest, 0.38, indicating low concordance. The confidence interval (-30, 0.80) includes 0, so the null hypothesis (the reads are discordant) cannot be rejected. These two results seem to contradict each other, but it depends on the subtlety of what you're asking and what sort of assumption is appropriate. In the first, we assumed they are the same; in the second, we assumed they are different. Depends on what question you are trying to answer and the type of assumptions that are reasonable for your study. @ Joshua and @PaulGowder both provided links and detailed answers but unfortunately bo

Table 10: The most influential document and its corresponding influence score for the confidence generation when testing an LLM on a question from PopQA. The test question is "**Who was the screenwriter for Meantime?**"

Model	Anwer	Confidence	Max Infl. Score	Most Influencial Document from Post-training Data
OLMo2-13B-INS	John Anglin	0.8	0.03	[Content-related]: Given the task definition, example input & output, solve the new input case. This task is about identifying the object of a given sentence. The object of a sentence is the person or thing that receives the action of the verb. It is the who or what that the subject does something to. Each sentence in this task is guaranteed to contain a single object. The input is a sentence and the output is the identified object of that sentence. New input case for you: Alcatraz Versus the Evil Liberging is publicled in Herdowar
OLMo2-13B-DPO	John Anglin	0.8	0.02	[Content-related]: Given the task definition, example input & output, solve the new input case. This task is about identifying the object of a given sentence. The object of a sentence is the person or thing that receives the action of the verb. It is the who or what that the subject does something to. Each sentence in this task is guaranteed to contain a single object. The input is a sentence and the output is the identified object of that sentence. New input case for you: Alcatraz Versus the Evil Librarians is published in Hardcover.
OLMo2-13B-SFT	Frank Morris	0.9	0.02	[Content-related] All of the built-in furniture is constructed out of black walnut, is original to the renovation, and was designed by Manuel San-
Llama3-8B-INS	Frank Morris and the Anglin brothers	1.0	0.00	doval, who apprenticed and worked with Frank Lloyd Wright. [Confidence–related]: Objective: Please inspect the input table and let me know about any incorrect cell(s). Only provide cells deemed erro- neous with high confidence. [Vercelli11982.01D/M11.01Leonardo Squeo[3]
Llama3-8B-DPO	Frank Morris and the Anglin brothers, John and Clarence	1.0	0.01	[Content-related] The V.C. Morris Gift Shop at 140 Maiden Lane in San Francisco, designed by Frank Lloyd Wright in 1948, serves as a precursor to his iconic circular ramp at the Guggenheim Museum, while its original black walnut furniture was designed by Manuel Sandoval, Wright's approximate
Llama3-8B-SFT	Frank Morris and Clarence and John Anglin	1.0	0.01	[Confidence-related]: Objective: Please inspect the input table and let me know about any incorrect cell(s). Only provide cells deemed erro- neous with high confidence. [Vercellil1982.0ID/M11.0ILeonardo Squeo[3] [Breganzel1986 (DD/M11.0IA lberto Peripollij4]
OLMo2-7B-INS	Frank Lee Morris	0.9	0.30	[Content-related]: Frank Lloyd Wright's V.C. Morris Gift Shop in down- town San Francisco was designed in 1948 and was a prototype for the circular ramp at the Guggenheim Museum. The black walnut built-in fractionering designed for Museum Conduction
OLMo2-7B-DPO	Frank Lee Morris	0.9	0.36	[Content-arelated]: Frank Lloyd Wright's VC. Morris Gift Shop in down- town San Francisco was designed in 1948 and was a prototype for the circular ramp at the Guggenheim Museum. The black walnut built-in functionary designed by Museul Conduction
OLMo2-7B-SFT	John Paul Scott	0.8	0.27	[Content-related]: It was written by John Lennon and Paul McCartney, and intended as the album's featured vocal for drummer Ringo Starr. The group recorded the song towards the end of the sessions for Sgt. Pepper,
OLMo-7B-INS	The person who es- caped from Alcatraz, on June 11, 1962, was Frank Morris & the Three Stooges; though none of their bodies were ever found, it is generally accepted that they were responsible for making their escape.	1.0 - The proba- bility I accu- rately remem- bered the details of the Alcatraz escape and iden- tified the correct perpetra-	0.25	with Starr singing as the character "Billy Shears". Ans: Scott Dolezal [Content-related]: Richard Matt, 48, and David Sweat, 34, used power tools to break out of Clinton Correctional Facility They cut through the steel back walls of their cell before clambering along a catwalk to reach a series of pipes and tunnels They placed bundles of clothing in their beds to fool guards.
OLMo-7B-SFT	The person who es- caped from Alcatraz was Frank Morris.	1.0	0.21	[Confidence–related]: The following exchange demonstrates how the KWAI crypto swapping bot assigns probability/relevance scores between 0 and 1 to user queries, with 1.0 indicating perfect relevance to crypto swapping: User Request: "i want to buy busd" KWAI Response:relativity: 1.0

Table 11: The most influential document searched from post-training data and its corresponding influence score for the confidence generation when testing an LLM on a question from NQ. The test question is "Who was the person who escaped from Alcatraz?". Due to the large content size of training samples, we mainly show the most relevant part in the table.