
Causally Quantifying the Effect of Test Set Contamination on Generative Benchmarks

Anonymous Author(s)

Affiliation
Address
email

Abstract

As large language models (LLMs) are pretrained on ever-expanding web-scale data, test set contamination has become a critical concern for accurately assessing the capabilities of LLMs. While significant research has quantified the amount and the impact of test set contamination on discriminative (i.e., scoring-based) benchmarks like multiple-choice question-answering, comparatively little research has studied the impact of test set contamination on generative (i.e., sampling-based) evaluations such as coding or mathematical problem solving. As the field shifts more towards generative evaluations, understanding what effect (if any) test set contamination has on generative evaluations becomes all the more important. To causally quantify the effect that test set contamination has on assessed capabilities, we pretrained language models, sweeping the number of replicas of benchmark test data in the pretraining corpora. We make four discoveries: (1) performance increases with contamination and model size, consistent with discriminative evaluations, (2) higher sampling temperature mitigates the effects of contamination, (3) longer solutions require more contamination to reach the same level of performance, and (4) generative performance is tightly coupled with test set memorization, but modulated by sampling temperature. As the field shifts to generative benchmarks to assess reasoning, our work reveals that factors like sampling temperature and solution length introduce novel complexities to data contamination, demanding a more sophisticated approach to model evaluation.

1 Introduction

Test set contamination – the inclusion of benchmarks in pretraining data – has emerged as a critical threat to the trustworthy evaluation of language models (Sainz et al., 2023; Schaeffer, 2023; Xu et al., 2024a; Deng et al., 2024a; Reuel et al., 2025). Evaluation aims to measure generalization on tasks the model has never seen, yet web-scale pretraining makes such contamination increasingly likely (Brown et al., 2020; Du et al., 2022; Wei et al., 2022; Chowdhery et al., 2022; Touvron et al., 2023).

Prior work has sought to quantify how training on benchmarks affects model evaluation scores through statistical and causal approaches. Statistical approaches aim to quantify the influence of test set contamination on evaluation performance by modifying the test set, for example, by reordering, rephrasing, or replicating benchmark problems, e.g., (Oren et al., 2023; Ni et al., 2025; Shi et al., 2024; Golchin & Surdeanu, 2023, 2024; Roberts et al., 2024; Wang et al., 2025; Zhang et al., 2024a). Causal approaches intentionally contaminate pretraining corpora to quantify how a particular dose of contamination translates into improved performance, e.g., (Magar & Schwartz, 2022; Jiang et al., 2024; Oren et al., 2023; Yao et al., 2024; Wang et al., 2025; Kocigit et al., 2025; Bordt et al., 2025).

While both are useful, increasing model capabilities and the advent of reasoning models (OpenAI et al., 2024; DeepMind et al., 2025; Xu et al., 2025) have shifted the field from discriminative

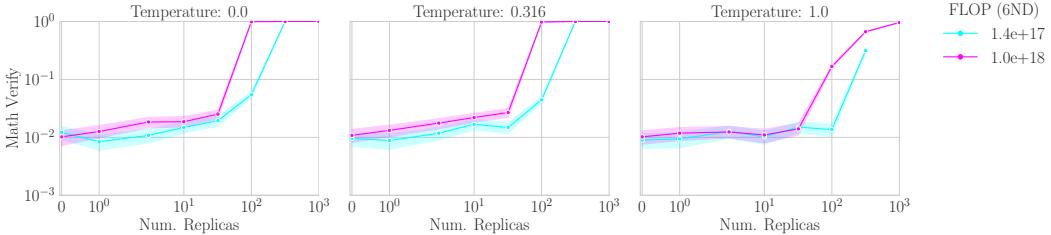


Figure 1: Causally Quantifying the Effect of Test Set Contamination on Generative Benchmarks. We pretrained language models of increasing size, while sweeping the number of generative benchmark replicas included in the pretraining corpora. We specifically chose MATH (Hendrycks et al., 2021b) as our generative benchmark of interest; Math Verify is the fraction of problems for which the model generates solutions verified to be mathematically equivalent to the benchmark’s boxed answers. Math Verify scores increase with the amount of contamination (i.e., the number of test set replicas) as well as with model scale. Consistent with discriminative evaluations, larger models require fewer replicas of the test set to reach ceiling performance. Left to Right: Higher sampling temperature mitigates the effects of contamination. Shaded regions represent 95% confidence intervals.

37 benchmarks to generative benchmarks, and research on test set contamination has lagged behind.
 38 Previous investigations on contamination focused on *discriminative* (i.e., scoring-based) benchmarks
 39 like classification or multiple-choice question-answering (MCQA), but newer evaluations prioritize
 40 *generative* (i.e., sampling-based) benchmarks like coding, mathematical problem solving or agentic
 41 tasks. For example, prior papers that intentionally contaminated pretraining corpora focused over-
 42 whelmingly on discriminative benchmarks: Magar & Schwartz (2022) used SST-2 (Socher et al.,
 43 2013) (classification). Jiang et al. (2024) used SST-2 (classification), MMLU (Hendrycks et al.,
 44 2021a) (MCQA), SQuAD (Rajpurkar et al., 2016) (MCQA), and CNN/Daily Mail (fill-in-the-middle)
 45 (Nallapati et al., 2016). Oren et al. (2023) used 7 MCQA benchmarks and 1 mathematical problem
 46 solving benchmark (GSM8K) (Cobbe et al., 2021), Yao et al. (2024) used 3 MCQA benchmarks
 47 while Bordt et al. (2025) used 7 MCQA benchmarks.

48 Whether test set contamination has the same effect on generative and discriminative evaluations is
 49 unclear *a priori*. Discriminative evaluations require the model to place higher probability mass on the
 50 correct choice than on a small number of alternative incorrect choices (Gao et al., 2024; Schaeffer
 51 et al., 2025b), and candidate choices are often only a couple of tokens long. In comparison, generative
 52 evaluations introduce more axes of choice into evaluation (i.e., how tokens are sampled and how the
 53 model is prompted) and require the model to produce tens to thousands of tokens, e.g., (Jimenez et al.,
 54 2024). While a sufficiently large model trained on a generative benchmark for an infinite number
 55 of epochs would almost certainly memorize its content (Carlini et al., 2023), it is currently unclear
 56 how much performance improves after a finite number of exposures; perhaps exponentially many
 57 repetitions are necessary to memorize such long sequences.

58 In this work, we quantify the effect that causally contaminating pretraining corpora with generative
 59 reasoning benchmarks has, focusing specifically on MATH (Hendrycks et al., 2021b). We find
 60 that: (1) performance increases with contamination and model size, consistent with discriminative
 61 evaluations, (2) higher sampling temperature might reduce the effects of contamination, (3) longer
 62 solutions require more contamination to reach the same level of performance, and (4) generative
 63 performance is tightly coupled with test set memorization, but is modulated by higher sampling
 64 temperatures.

65 2 Experimental Setup

66 **Benchmark** For our generative benchmark of interest, we chose MATH (Hendrycks et al., 2021b);
 67 its test set contains approximately 1.5M tokens.

68 **Pretraining** We pretrained Qwen 3 (Yang et al., 2025a) transformer-based (Vaswani et al., 2017)
 69 causal language models from initialization. In this work, we consider two model sizes, 34M and
 70 93M parameters, as our 344M and 1.44B parameter models are still pretraining at the time of this

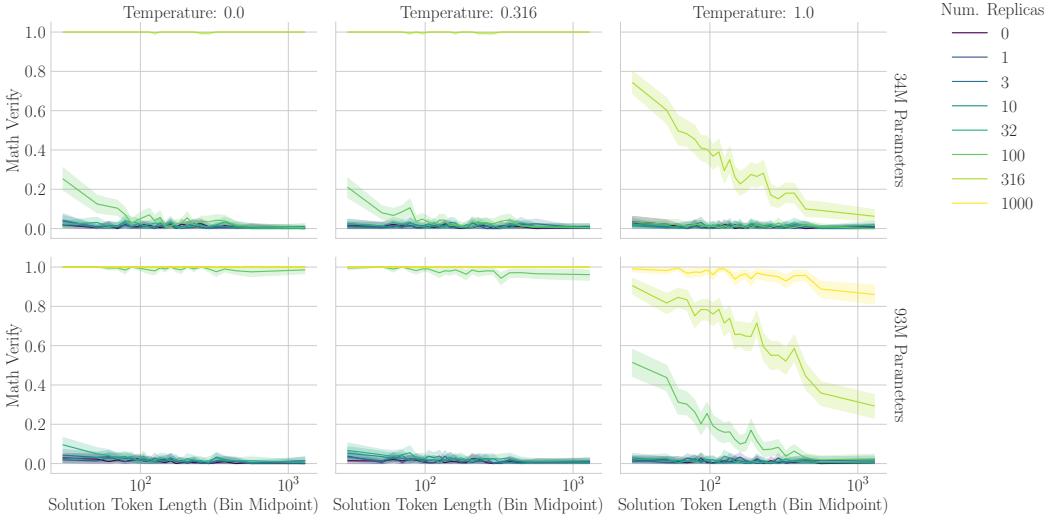


Figure 2: Model Performance Declines with Increasing Solution Token Length. Math Verify scores depend primarily on model size and the number of benchmark replicas. However, scores clearly decrease as a function of the number of tokens in the solution, and higher temperature amplifies this effect (right). Shaded regions represent 95% confidence intervals.

71 submission. For each model size, we created multiple pretraining corpora by concatenating FineWeb-
 72 Edu (Penedo et al., 2024) with a different number of replicas of the benchmark test set from 0
 73 (uncontaminated) through 1, 3, 10, 32, 100, 316, 1000. For our benchmark of interest, we chose the
 74 MATH (Hendrycks et al., 2021b) test set, containing approximately 1.5M tokens. Each model was
 75 pretrained for 20 tokens per parameter, following compute-optimal scaling (Hoffmann et al., 2022).
 76 Pretraining compute was calculated using the common approximation $C \approx 6 N D$ (Kaplan et al.,
 77 2020; Porian et al., 2024), where N is the number of parameters and D is the number of tokens.

78 **Evaluating** We evaluated our models using EleutherAI’s Language Model Evaluation Harness
 79 (Gao et al., 2024), which reports the “math verify” score: the fraction of problems for which the
 80 model generates solutions that are verified to be mathematically equivalent to the benchmark’s boxed
 81 answers. Along the way, we discovered an error with how the Harness computes Math Verify scores
 82 on MATH, and worked with its maintainers to fix the error; this suggests to us that any research
 83 reporting MATH scores from the past 1+ years may have egregiously incorrect values. We used
 84 basic (i.e., temperature-only) sampling (Schaeffer et al., 2025a) and report scores for three different
 85 sampling temperatures: 0 (greedy), 0.316 and 1.00.

86 3 Results

87 In this section, we report how model performance (i.e., Math Verify score) changes as a function of
 88 model size, number of test set replicas, and generative evaluation sampling temperature.

89 **Finding #1: Performance Increases with Contamination and Model Size** Consistent with
 90 discriminative evaluations, we find that increasing the number of replicas of the benchmark in the
 91 pretraining corpus increases Math Verify scores, as does increasing the model size (Fig 1). We
 92 observe a non-linear relationship between the number of test set replicas and model performance: For
 93 low levels of contamination (≤ 10 replicas), the impact on performance is minimal, with Math Verify
 94 scores remaining close to the baseline (0 replicas) performance; at around 100 replicas, a notable
 95 inflection point occurs, where performance sharply increases across all temperatures; at the highest
 96 level of contamination (316 replicas), the model achieves near-perfect performance, particularly at
 97 lower temperatures. Although we currently have only two model sizes, the trend suggests that while
 98 minor contamination may not significantly affect evaluation of tiny models, minor contamination of
 99 larger models can lead to severe overestimation of model capabilities.

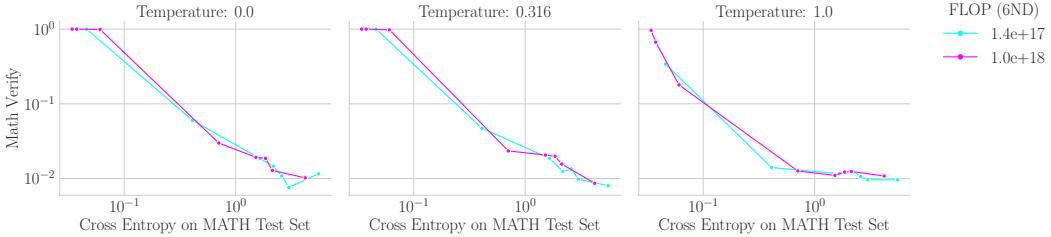


Figure 3: Generative Performance is Tightly Coupled with Test Set Memorization, but is Modulated by Sampling Temperature. As the number of test set replicas in the pretraining data increases, the cross-entropy loss decreases, leading to a sharp increase in the Math Verify score. This indicates that improved generative performance is tightly coupled with the model’s memorization of the test set, as measured by cross-entropy. Notably, a higher sampling temperature increases the curvature of this relationship (right), resulting in a less abrupt performance increase as cross-entropy decreases. Shaded regions represent 95% confidence intervals.

100 **Finding #2: Higher Sampling Temperature Might Mitigate the Effects of Contamination** Math
 101 Verify Scores do not change substantially between zero-temperature and low temperature sampling
 102 (0.316) (Fig. 1, left and center). At higher temperature (1.0), Math Verify scores are reduced relative to
 103 lower temperatures (0 and 0.316). Interestingly, higher temperature sampling reduces the gap between
 104 highly contaminated models and uncontaminated models, suggesting that higher temperatures might
 105 mitigate the effect of contamination. This relationship between temperature and extraction is not
 106 always monotonic; other work has found that the optimal temperature for extracting memorized text
 107 can vary and is co-dependent on other details like the sampling algorithm (Hayes et al., 2025).

108 **Finding #3: Longer Solutions Require More Contamination To Reach the Same Performance**
 109 To understand how the length of the solution affects model performance, we bin problems based on
 110 solution token length and compute the average Math Verify score per bin. While model performance
 111 depends primarily on model size and number of test set replicas, scores decrease with the token
 112 length of the solution (Fig 2). Higher temperature additionally amplifies this effect. We currently do
 113 not have sufficient data to mathematically describe the relationships, but we intend to pursue this in
 114 future work once our larger models finish pretraining. Prior work on memorization corroborates this
 115 claim, as Jiang et al. (2025) and Lu et al. (2024) find that longer memorized sequences are the higher
 116 repeated ones, although to the best of our knowledge, this has not been studied causally in prior work.

117 **Finding #4: Generative Performance is Tightly Coupled with Test Set Memorization, but Is**
 118 **Modulated by Sampling Temperature** We find a strong negative correlation between a model’s
 119 Math Verify score and its cross-entropy loss on the MATH test set (Fig. 3). As we increase the number
 120 of benchmark replicas in the pretraining data, the model’s cross-entropy on the test set decreases,
 121 which in turn is associated with a sharp, non-linear increase in its generative performance. This
 122 tight coupling suggests that the performance gains are a direct consequence of memorization, with
 123 cross-entropy on the test set serving as a quantitative proxy for this effect. This effect is moderated
 124 by sampling temperature; higher temperatures increase the curvature of the relationship, resulting in
 125 a less abrupt performance increase as the model better memorizes the test set (Fig. 3, right).

126 4 Discussion

127 Our findings deliver a crucial caution as the field increasingly turns to generative benchmarks to
 128 evaluate advanced reasoning. We causally demonstrate that performance gains on these benchmarks
 129 are not necessarily evidence of improved reasoning, but are tightly coupled with test set memorization.
 130 This relationship, however, is not straightforward. Unlike in discriminative tasks, the effects of
 131 contamination are modulated by new factors unique to generative evaluation: higher sampling
 132 temperatures can mitigate or mask the impact of contamination, while longer solutions require more
 133 contaminating repetitions to achieve the same inflated performance. These complexities introduce a
 134 significant risk of severely overestimating model capabilities.

135 **Related Work** Due to space limitations, we defer related work to Appendix A.

136 **References**

- 137 Ben Adlam and Jeffrey Pennington. Understanding double descent requires a fine-grained bias-
138 variance decomposition. *Advances in neural information processing systems*, 33:11022–11032,
139 2020.
- 140 Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of general-
141 ization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- 142 Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning
143 practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*,
144 116(32):15849–15854, 2019.
- 145 Stella Biderman, Usvsn Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivan-
146 shu Purohit, and Edward Raff. Emergent and predictable memorization in large language models.
147 *Advances in Neural Information Processing Systems*, 36:28072–28090, 2023.
- 148 Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in
149 kernel regression and wide neural networks. In *International Conference on Machine Learning*,
150 pp. 1024–1034. PMLR, 2020.
- 151 Sebastian Bordt, Suraj Srinivas, Valentyn Boreiko, and Ulrike von Luxburg. How much can we forget
152 about data contamination? In *Forty-second International Conference on Machine Learning*, 2025.
- 153 Dillon Bowen, Brendan Murphy, Will Cai, David Khachaturov, Adam Gleave, and Kellin Pelrine.
154 Scaling trends for data poisoning in llms, 2025. URL <https://arxiv.org/abs/2408.02946>.
- 155 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
156 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
157 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 158 Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine
159 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data
160 from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pp.
161 2633–2650, 2021.
- 162 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan
163 Zhang. Quantifying memorization across neural language models. In *The Eleventh International
164 Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=TatRHT_1cK.
- 166 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
167 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,
168 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam
169 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James
170 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Lev-
171 skaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin
172 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph,
173 Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M.
174 Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon
175 Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark
176 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean,
177 Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2022. URL
178 <https://arxiv.org/abs/2204.02311>.
- 179 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
180 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
181 Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- 183 Debeshee Das, Jie Zhang, and Florian Tramèr. Blind baselines beat membership inference attacks for
184 foundation models. *arXiv preprint arXiv:2406.16201*, 2024.

185 Google DeepMind, Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen
186 Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris,
187 Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson,
188 Idan Szpektor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri,
189 Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornraphop
190 Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilai Deutel,
191 Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare,
192 Helen Miller, Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero
193 Rissa, Dan Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves,
194 Cindy Wang, Nitesh Bharadwaj Gundavarapu, Ilia Shumailov, Bo Wang, Mantas Pajarskas, Joe
195 Heyward, Martin Nikoltchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik,
196 Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kavukcuoglu,
197 Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov,
198 Bo-Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas
199 Leichner, Haichuan Yang, Zelda Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin
200 Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios
201 Vytniotis, Jieru Mei, Mu Cai, Mohammed Badawi, Corey Fry, Ale Hartman, Daniel Zheng, Eric
202 Jia, James Keeling, Annie Louis, Ying Chen, Efren Robles, Wei-Chih Hung, Howard Zhou, Nikita
203 Saxena, Sonam Goenka, Olivia Ma, Zach Fisher, Mor Hazan Taege, Emily Graves, David Steiner,
204 Yujia Li, Sarah Nguyen, Rahul Sukthankar, Joe Stanton, Ali Eslami, Gloria Shen, Berkin Akin,
205 Alexey Guseynov, Yiqian Zhou, Jean-Baptiste Alayrac, Armand Joulin, Efrat Farkash, Ashish
206 Thapliyal, Stephen Roller, Noam Shazeer, Todor Davchev, Terry Koo, Hannah Forbes-Pollard,
207 Kartik Audhkhasi, Greg Farquhar, Adi Mayrav Gilady, Maggie Song, John Aslanides, Piermaria
208 Mendolicchio, Alicia Parrish, John Blitzer, Pramod Gupta, Xiaoen Ju, Xiaochen Yang, Puranjay
209 Datta, Andrea Tacchetti, Sanket Vaibhav Mehta, Gregory Dibb, Shubham Gupta, Federico Piccinini,
210 Raia Hadsell, Sujee Rajayogam, Jiepu Jiang, Patrick Griffin, Patrik Sundberg, Jamie Hayes, Alexey
211 Frolov, Tian Xie, Adam Zhang, Kingshuk Dasgupta, Uday Kalra, Lior Shani, Klaus Macherey,
212 Tzu-Kuo Huang, Liam MacDermid, Karthik Duddu, Paulo Zacchello, Zi Yang, Jessica Lo, Kai
213 Hui, Matej Kastelic, Derek Gasaway, Qijun Tan, Summer Yue, Pablo Barrio, John Wieting, Weel
214 Yang, Andrew Nystrom, Solomon Demmessie, Anselm Levskaya, Fabio Viola, Chetan Tekur,
215 Greg Billock, George Necula, Mandar Joshi, Rylan Schaeffer, Swachhand Lokhande, Christina
216 Sorokin, Pradeep Shenoy, Mia Chen, Mark Collier, Hongji Li, Taylor Bos, Nevan Wichers,
217 Sun Jae Lee, Angéline Pouget, Santhosh Thangaraj, Kyriakos Axiotis, Phil Crone, Rachel Sterneck,
218 Nikolai Chinaev, Victoria Krakovna, Oleksandr Ferludin, Ian Gemp, Stephanie Winkler, Dan
219 Goldberg, Ivan Korotkov, Kefan Xiao, Malika Mehrotra, Sandeep Mariserla, Vihari Piratla, Terry
220 Thurk, Khiem Pham, Hongxu Ma, Alexandre Senges, Ravi Kumar, Clemens Meyer, Ellie Talius,
221 Nuo Wang Pierse, Ballie Sandhu, Horia Toma, Kuo Lin, Swaroop Nath, Tom Stone, Dorsa Sadigh,
222 Nikita Gupta, Arthur Guez, Avi Singh, Matt Thomas, Tom Duerig, Yuan Gong, Richard Tanburn,
223 Lydia Lihui Zhang, Phuong Dao, Mohamed Hammad, Sirui Xie, Shruti Rijhwani, Ben Murdoch,
224 Duhyeon Kim, Will Thompson, Heng-Tze Cheng, Daniel Sohn, Pablo Sprechmann, Qiantong
225 Xu, Srinivas Tadepalli, Peter Young, Ye Zhang, Hansa Srinivasan, Miranda Aperghis, Aditya
226 Ayyar, Hen Fitoussi, Ryan Burnell, David Madras, Mike Dusenberry, Xi Xiong, Tayo Oguntebi,
227 Ben Albrecht, Jörg Bornschein, Jovana Mitrović, Mason Dimarco, Bhargav Kanagal Shamanna,
228 Premal Shah, Eren Sezener, Shyam Upadhyay, Dave Lacey, Craig Schiff, Sebastien Baur, Sanjay
229 Ganapathy, Eva Schnider, Mateo Wirth, Connor Schenck, Andrey Simanovsky, Yi-Xuan Tan,
230 Philipp Fränken, Dennis Duan, Bharath Mankalale, Nikhil Dhawan, Kevin Sequeira, Zichuan
231 Wei, Shivanker Goel, Caglar Unlu, Yukun Zhu, Haitian Sun, Ananth Balashankar, Kurt Shuster,
232 Megh Umekar, Mahmoud Alnahlawi, Aäron van den Oord, Kelly Chen, Yuexiang Zhai, Zihang
233 Dai, Kuang-Huei Lee, Eric Doi, Lukas Zilka, Rohith Vallu, Disha Shrivastava, Jason Lee, Hisham
234 Husain, Honglei Zhuang, Vincent Cohen-Addad, Jarred Barber, James Atwood, Adam Sadovsky,
235 Quentin Wellens, Steven Hand, Arunkumar Rajendran, Aybuke Turker, CJ Carey, Yuanzhong Xu,
236 Hagen Soltau, Zefei Li, Xinying Song, Conglong Li, Jurii Kemaev, Sasha Brown, Andrea Burns,
237 Viorica Patrascu, Piotr Stanczyk, Renga Aravamudhan, Mathieu Blondel, Hila Noga, Lorenzo
238 Blanco, Will Song, Michael Isard, Mandar Sharma, Reid Hayes, Dalia El Badawy, Avery Lamp,
239 Itay Laish, Olga Kozlova, Kelvin Chan, Sahil Singla, Srinivas Sunkara, Mayank Upadhyay, Chang
240 Liu, Aijun Bai, Jarek Wilkiewicz, Martin Zlocha, Jeremiah Liu, Zhuowan Li, Haiguang Li, Omer
241 Barak, Ganna Raboshchuk, Jiho Choi, Fangyu Liu, Erik Jue, Mohit Sharma, Andreea Marzoca,
242 Robert Busa-Fekete, Anna Korsun, Andre Elisseeff, Zhe Shen, Sara Mc Carthy, Kay Lamrigts,
243 Anahita Hosseini, Hanzhao Lin, Charlie Chen, Fan Yang, Kushal Chauhan, Mark Omernick,

244 Dawei Jia, Karina Zainullina, Demis Hassabis, Danny Vainstein, Ehsan Amid, Xiang Zhou, Ronny
245 Votel, Eszter Vértes, Xinjian Li, Zongwei Zhou, Angeliki Lazaridou, Brendan McMahan, Arjun
246 Narayanan, Hubert Soyer, Sujoy Basu, Kayi Lee, Bryan Perozzi, Qin Cao, Leonard Berrada, Rahul
247 Arya, Ke Chen, Katrina, Xu, Matthias Lochbrunner, Alex Hofer, Sahand Sharifzadeh, Renjie
248 Wu, Sally Goldman, Pranjal Awasthi, Xuezhi Wang, Yan Wu, Claire Sha, Biao Zhang, Maciej
249 Mikuła, Filippo Graziano, Siobhan McIoughlin, Irene Giannoumis, Youhei Namiki, Chase Malik,
250 Carey Radebaugh, Jamie Hall, Ramiro Leal-Cavazos, Jianmin Chen, Vikas Sindhwani, David Kao,
251 David Greene, Jordan Griffith, Chris Welty, Ceslee Montgomery, Toshihiro Yoshino, Liangzhe
252 Yuan, Noah Goodman, Assaf Hurwitz Michaely, Kevin Lee, KP Sawhney, Wei Chen, Zheng
253 Zheng, Megan Shum, Nikolay Savinov, Etienne Pot, Alex Pak, Morteza Zadimoghaddam, Sijal
254 Bhatnagar, Yoad Lewenberg, Blair Kutzman, Ji Liu, Lesley Katzen, Jeremy Selier, Josip Djolonga,
255 Dmitry Lepikhin, Kelvin Xu, Jacky Liang, Jiewen Tan, Benoit Schillings, Muge Ersoy, Pete
256 Blois, Bernd Bandemer, Abhimanyu Singh, Sergei Lebedev, Pankaj Joshi, Adam R. Brown, Evan
257 Palmer, Shreya Pathak, Komal Jalan, Fedir Zubach, Shuba Lall, Randall Parker, Alok Gunjan,
258 Sergey Rogulenko, Sumit Sanghai, Zhaoqi Leng, Zoltan Egyed, Shixin Li, Maria Ivanova, Kostas
259 Andriopoulos, Jin Xie, Elan Rosenfeld, Auriel Wright, Ankur Sharma, Xinyang Geng, Yicheng
260 Wang, Sam Kwei, Renke Pan, Yujing Zhang, Gabby Wang, Xi Liu, Chak Yeung, Elizabeth
261 Cole, Aviv Rosenberg, Zhen Yang, Phil Chen, George Polovets, Pranav Nair, Rohun Saxena,
262 Josh Smith, Shuo yin Chang, Aroma Mahendru, Svetlana Grant, Anand Iyer, Irene Cai, Jed
263 McGiffin, Jiaming Shen, Alanna Walton, Antonious Girgis, Oliver Woodman, Rosemary Ke, Mike
264 Kwong, Louis Rouillard, Jimmeng Rao, Zhihao Li, Yuntao Xu, Flavien Prost, Chi Zou, Ziwei Ji,
265 Alberto Magni, Tyler Liechty, Dan A. Calian, Deepak Ramachandran, Igor Krivokon, Hui Huang,
266 Terry Chen, Anja Hauth, Anastasija Ilić, Weijuan Xi, Hyeontaek Lim, Vlad-Doru Ion, Pooya
267 Moradi, Metin Toksoz-Exley, Kalesha Bullard, Miltos Allamanis, Xiaomeng Yang, Sophie Wang,
268 Zhi Hong, Anita Gergely, Cheng Li, Bhavishya Mittal, Vitaly Kovalev, Victor Ungureanu, Jane
269 Labanowski, Jan Wassenberg, Nicolas Lacasse, Geoffrey Cideron, Petar Dević, Annie Marsden,
270 Lynn Nguyen, Michael Fink, Yin Zhong, Tatsuya Kiyono, Desi Ivanov, Sally Ma, Max Bain,
271 Kiran Yalasangi, Jennifer She, Anastasia Petrushkina, Mayank Lunayach, Carla Bromberg, Sarah
272 Hodkinson, Vilobh Meshram, Daniel Vlasic, Austin Kyker, Steve Xu, Jeff Stanway, Zuguang Yang,
273 Kai Zhao, Matthew Tung, Seth Odoom, Yasuhisa Fujii, Justin Gilmer, Eunyoung Kim, Felix Halim,
274 Quoc Le, Bernd Bohnet, Seliem El-Sayed, Behnam Neyshabur, Malcolm Reynolds, Dean Reich,
275 Yang Xu, Erica Moreira, Anuj Sharma, Zeyu Liu, Mohammad Javad Hosseini, Naina Raisinghani,
276 Yi Su, Ni Lao, Daniel Formoso, Marco Gelmi, Almog Gueta, Tapomay Dey, Elena Gribovskaya,
277 Domagoj Ćevid, Sidharth Mudgal, Garrett Bingham, Jianling Wang, Anurag Kumar, Alex Cullum,
278 Feng Han, Konstantinos Bousmalis, Diego Cedillo, Grace Chu, Vladimir Magay, Paul Michel,
279 Ester Hlavnova, Daniele Calandriello, Setareh Ariafar, Kaisheng Yao, Vikash Sehwag, Arpi Vezer,
280 Agustin Dal Lago, Zhenkai Zhu, Paul Kishan Rubenstein, Allen Porter, Anirudh Baddepudi, Oriana
281 Riva, Mihai Dorin Istiń, Chih-Kuan Yeh, Zhi Li, Andrew Howard, Nilpa Jha, Jeremy Chen, Raoul
282 de Liedekerke, Zafarali Ahmed, Mikel Rodriguez, Tanuj Bhatia, Bangju Wang, Ali Elqursh, David
283 Klinghoffer, Peter Chen, Pushmeet Kohli, Te I, Weiyang Zhang, Zack Nado, Jilin Chen, Maxwell
284 Chen, George Zhang, Ayush Singh, Adam Hillier, Federico Lebron, Yiqing Tao, Ting Liu, Gabriel
285 Dulac-Arnold, Jingwei Zhang, Shashi Narayan, Buhuang Liu, Orhan Firat, Abhishek Bhowmick,
286 Bingyuan Liu, Hao Zhang, Zizhao Zhang, Georges Rotival, Nathan Howard, Anu Sinha, Alexander
287 Grushetsky, Benjamin Beyret, Keerthana Gopalakrishnan, James Zhao, Kyle He, Szabolcs Payrits,
288 Zaid Nabulsi, Zhaoyi Zhang, Weijie Chen, Edward Lee, Nova Fallen, Sreenivas Gollapudi, Aurick
289 Zhou, Filip Pavetić, Thomas Köppe, Shiyu Huang, Rama Pasumarthi, Nick Fernando, Felix
290 Fischer, Daria Ćurko, Yang Gao, James Svensson, Austin Stone, Haroon Qureshi, Abhishek
291 Sinha, Apoorv Kulshreshtha, Martin Matysiak, Jieming Mao, Carl Saroufim, Aleksandra Faust,
292 Qingnan Duan, Gil Fidel, Kaan Katircioğlu, Raphaël Lopez Kaufman, Dhruv Shah, Weize Kong,
293 Abhishek Bapna, Gellért Weisz, Emma Dunleavy, Praneet Dutta, Tianqi Liu, Rahma Chaabouni,
294 Carolina Parada, Marcus Wu, Alexandra Belias, Alessandro Bissacco, Stanislav Fort, Li Xiao,
295 Fantine Huot, Chris Knutson, Yochai Blau, Gang Li, Jennifer Prendki, Juliette Love, Yinlam
296 Chow, Pichi Charoenpanit, Hidetoshi Shimokawa, Vincent Coriou, Karol Gregor, Tomas Izo, Arjun
297 Akula, Mario Pinto, Chris Hahn, Dominik Paulus, Jiaxian Guo, Neha Sharma, Cho-Jui Hsieh,
298 Adaeze Chukwuka, Kazuma Hashimoto, Nathalie Rauschmayr, Ling Wu, Christof Angermueller,
299 Yulong Wang, Sebastian Gerlach, Michael Pliskin, Daniil Mirylenka, Min Ma, Lexi Baugher,
300 Bryan Gale, Shaan Bijwadia, Nemanja Rakićević, David Wood, Jane Park, Chung-Ching Chang,
301 Babi Seal, Chris Tar, Kacper Krasowiak, Yiwen Song, Georgi Stephanov, Gary Wang, Marcello
302 Maggioni, Stein Xudong Lin, Felix Wu, Shachi Paul, Zixuan Jiang, Shubham Agrawal, Bilal Piot,

303 Alex Feng, Cheolmin Kim, Tulsee Doshi, Jonathan Lai, Chuqiao, Xu, Sharad Vikram, Ciprian
304 Chelba, Sebastian Krause, Vincent Zhuang, Jack Rae, Timo Denk, Adrian Collister, Lotte Weerts,
305 Xianghong Luo, Yifeng Lu, Håvard Garnes, Nitish Gupta, Terry Spitz, Avinatan Hassidim, Lihao
306 Liang, Izhak Shafran, Peter Humphreys, Kenny Vassigh, Phil Wallis, Virat Shejwalkar, Nicolas
307 Perez-Nieves, Rachel Hornung, Melissa Tan, Beka Westberg, Andy Ly, Richard Zhang, Brian
308 Farris, Jongbin Park, Alec Kosik, Zeynep Cankara, Andrii Maksai, Yunhan Xu, Albin Cassirer,
309 Sergi Caelles, Abbas Abdolmaleki, Mencher Chiang, Alex Fabrikant, Shravya Shetty, Luheng
310 He, Mai Giménez, Hadi Hashemi, Sheena Panthapacakel, Yana Kulizhskaya, Salil Deshmukh,
311 Daniele Pighin, Robin Alazard, Disha Jindal, Seb Noury, Pradeep Kumar S, Siyang Qin, Xerxes
312 Dotiwalla, Stephen Spencer, Mohammad Babaeizadeh, Blake JianHang Chen, Vaibhav Mehta,
313 Jennie Lees, Andrew Leach, Penporn Koanantakool, Ilia Akolzin, Ramona Comanescu, Junwhan
314 Ahn, Alexey Svyatkovskiy, Basil Mustafa, David D'Ambrosio, Shiva Mohan Reddy Garlapati,
315 Pascal Lamblin, Alekh Agarwal, Shuang Song, Pier Giuseppe Sessa, Pauline Coquinot, John
316 Maggs, Hussain Masoom, Divya Pitta, Yaqing Wang, Patrick Morris-Suzuki, Billy Porter, Johnson
317 Jia, Jeffrey Dudek, Raghavender R, Cosmin Paduraru, Alan Ansell, Tolga Bolukbasi, Tony Lu,
318 Ramya Ganeshan, Zi Wang, Henry Griffiths, Rodrigo Benenson, Yifan He, James Swirhun, George
319 Papamakarios, Aditya Chawla, Kuntal Sengupta, Yan Wang, Vedrana Milutinovic, Igor Mordatch,
320 Zhipeng Jia, Jamie Smith, Will Ng, Shitij Nigam, Matt Young, Eugen Vušak, Blake Hechtman,
321 Sheela Goenka, Avital Zipori, Kareem Ayoub, Ashok Popat, Trilok Acharya, Luo Yu, Dawn
322 Bloxwich, Hugo Song, Paul Roit, Haiqiong Li, Aviel Boag, Nigamaa Nayakanti, Bilva Chandra,
323 Tianli Ding, Aahil Mehta, Cath Hope, Jiageng Zhang, Idan Heimlich Shtacher, Kartikeya Badola,
324 Ryo Nakashima, Andrei Sozanschi, Iulia Comşa, Ante Žužul, Emily Caveness, Julian Odell,
325 Matthew Watson, Dario de Cesare, Phillip Lippe, Derek Lockhart, Siddharth Verma, Huizhong
326 Chen, Sean Sun, Lin Zhuo, Aditya Shah, Prakhar Gupta, Alex Muzio, Ning Niu, Amir Zait,
327 Abhinav Singh, Meenu Gaba, Fan Ye, Prajit Ramachandran, Mohammad Saleh, Raluca Ada Popa,
328 Ayush Dubey, Frederick Liu, Sara Javanmardi, Mark Epstein, Ross Hemsley, Richard Green,
329 Nishant Ranka, Eden Cohen, Chuyuan Kelly Fu, Sanjay Ghemawat, Jed Borovik, James Martens,
330 Anthony Chen, Pranav Shyam, André Susano Pinto, Ming-Hsuan Yang, Alexandru Tifrea, David
331 Du, Boqing Gong, Ayushi Agarwal, Seungyeon Kim, Christian Frank, Saloni Shah, Xiaodan Song,
332 Zhiwei Deng, Ales Mikhalap, Kleopatra Chatziprimou, Timothy Chung, Toni Creswell, Susan
333 Zhang, Yennie Jun, Carl Lebsack, Will Truong, Slavica Andačić, Itay Yona, Marco Fornoni, Rong
334 Rong, Serge Toropov, Afzal Shama Soudagar, Andrew Audibert, Salah Zaiem, Zaheer Abbas,
335 Andrei Rusu, Sahitya Potluri, Shitao Weng, Anastasios Kementsietsidis, Anton Tsitsulin, Daiyi
336 Peng, Natalie Ha, Sanil Jain, Tejasji Latkar, Simeon Ivanov, Cory McLean, Anirudh GP, Rajesh
337 Venkataraman, Canoeo Liu, Dilip Krishnan, Joel D'sa, Roey Yogev, Paul Collins, Benjamin Lee,
338 Lewis Ho, Carl Doersch, Gal Yona, Shawn Gao, Felipe Tiengo Ferreira, Adnan Ozturel, Hannah
339 Muckenheim, Ce Zheng, Gargi Balasubramaniam, Mudit Bansal, George van den Driessche, Sivan
340 Eiger, Salem Haykal, Vedant Misra, Abhimanyu Goyal, Danilo Martins, Gary Leung, Jonas
341 Valfridsson, Four Flynn, Will Bishop, Chenxi Pang, Yoni Halpern, Honglin Yu, Lawrence Moore,
342 Yuvein, Zhu, Sridhar Thiagarajan, Yoel Drori, Zhisheng Xiao, Lucio Dery, Rolf Jagerman, Jing
343 Lu, Eric Ge, Vaibhav Aggarwal, Arjun Khare, Vinh Tran, Oded Elyada, Ferran Alet, James Rubin,
344 Ian Chou, David Tian, Libin Bai, Lawrence Chan, Lukasz Lew, Karolis Misiunas, Taylan Bilal,
345 Aniket Ray, Sindhu Raghuram, Alex Castro-Ros, Viral Carpenter, CJ Zheng, Michael Kilgore,
346 Josef Broder, Emily Xue, Praveen Kallakuri, Dheeru Dua, Nancy Yuen, Steve Chien, John Schultz,
347 Saurabh Agrawal, Reut Tsarfaty, Jingcao Hu, Ajay Kannan, Dror Marcus, Nisarg Kothari, Baochen
348 Sun, Ben Horn, Matko Bošnjak, Ferjad Naeem, Dean Hirsch, Lewis Chiang, Boya Fang, Jie Han,
349 Qifei Wang, Ben Hora, Antoine He, Mario Lučić, Beer Changpinyo, Anshuman Tripathi, John
350 Youssef, Chester Kwak, Philippe Schlettner, Cat Graves, Rémi Leblond, Wenjun Zeng, Anders
351 Andreassen, Gabriel Rasskin, Yue Song, Eddie Cao, Junhyuk Oh, Matt Hoffman, Wojtek Skut,
352 Yichi Zhang, Jon Stritar, Xingyu Cai, Saarthak Khanna, Kathie Wang, Shriya Sharma, Christian
353 Reisswig, Younghoon Jun, Aman Prasad, Tatiana Sholokhova, Preeti Singh, Adi Gerzi Rosenthal,
354 Anian Ruoss, Françoise Beaufays, Sean Kirmani, Dongkai Chen, Johan Schalkwyk, Jonathan
355 Herzig, Been Kim, Josh Jacob, Damien Vincent, Adrian N Reyes, Ivana Balazevic, Léonard
356 Hussenot, Jon Schneider, Parker Barnes, Luis Castro, Spandana Raj Babbula, Simon Green,
357 Serkan Cabi, Nico Duduta, Danny Driess, Rich Galt, Noam Velan, Junjie Wang, Hongyang Jiao,
358 Matthew Mauger, Du Phan, Miteyan Patel, Vlado Galić, Jerry Chang, Eyal Marcus, Matt Harvey,
359 Julian Salazar, Elahe Dabir, Suraj Satishkumar Sheth, Amol Mandhane, Hanie Sedghi, Jeremiah
360 Willcock, Amir Zandieh, Shruthi Prabhakara, Aida Amini, Antoine Miech, Victor Stone, Massimo
361 Nicosia, Paul Niemczyk, Ying Xiao, Lucy Kim, Sławek Kwasiborski, Vikas Verma, Ada Maksutaj

362 Oflazer, Christoph Hirnschall, Peter Sung, Lu Liu, Richard Everett, Michiel Bakker, Ágoston
363 Weisz, Yufei Wang, Vivek Sampathkumar, Uri Shaham, Bibo Xu, Yasemin Altun, Mingqiu Wang,
364 Takaaki Saeki, Guanjie Chen, Emanuel Taropa, Shanthal Vasanth, Sophia Austin, Lu Huang,
365 Goran Petrovic, Qingyun Dou, Daniel Golovin, Grigory Rozhdestvenskiy, Allie Culp, Will Wu,
366 Motoki Sano, Divya Jain, Julia Proskurnia, Sébastien Cevey, Alejandro Cruzado Ruiz, Piyush
367 Patil, Mahdi Mirzazadeh, Eric Ni, Javier Snaider, Lijie Fan, Alexandre Fréchette, AJ Piergiovanni,
368 Shariq Iqbal, Kenton Lee, Claudio Fantacci, Jinwei Xing, Lisa Wang, Alex Irpan, David Raposo,
369 Yi Luan, Zhuoyuan Chen, Harish Ganapathy, Kevin Hui, Jiazhong Nie, Isabelle Guyon, Heming
370 Ge, Roopali Vij, Hui Zheng, Dayeong Lee, Alfonso Castaño, Khuslen Baatarsukh, Gabriel
371 Ibagon, Alexandra Chronopoulou, Nicholas FitzGerald, Shashank Viswanadha, Safeen Huda,
372 Rivka Moroshko, Georgi Stoyanov, Prateek Kolhar, Alain Vaucher, Ishaan Watts, Adhi Kuncoro,
373 Henryk Michalewski, Satish Kambala, Bat-Orgil Batsaikhan, Alek Andreev, Irina Jurenka, Maigo
374 Le, Qihang Chen, Wael Al Jishi, Sarah Chakera, Zhe Chen, Aditya Kini, Vikas Yadav, Aditya
375 Siddhant, Ilia Labzovsky, Balaji Lakshminarayanan, Carrie Grimes Bostock, Pankil Botadra,
376 Ankesh Anand, Colton Bishop, Sam Conway-Rahman, Mohit Agarwal, Yani Donchev, Achintya
377 Singhal, Félix de Chaumont Quirity, Natalia Ponomareva, Nishant Agrawal, Bin Ni, Kalpesh
378 Krishna, Masha Samsikova, John Karro, Yilun Du, Tamara von Glehn, Caden Lu, Christopher A.
379 Choquette-Choo, Zhen Qin, Tingnan Zhang, Sicheng Li, Divya Tyam, Swaroop Mishra, Wing
380 Lowe, Colin Ji, Weiyi Wang, Manaal Faruqui, Ambrose Slone, Valentin Dalibard, Arunachalam
381 Narayanaswamy, John Lambert, Pierre-Antoine Manzagol, Dan Karliner, Andrew Bolt, Ivan
382 Lobov, Aditya Kusupati, Chang Ye, Xuan Yang, Heiga Zen, Nelson George, Mukul Bhutani,
383 Olivier Lacombe, Robert Riachi, Gagan Bansal, Rachel Soh, Yue Gao, Yang Yu, Adams Yu,
384 Emily Nottage, Tania Rojas-Esponda, James Noraky, Manish Gupta, Ragha Kotikalapudi, Jichuan
385 Chang, Sanja Deur, Dan Graur, Alex Mossin, Erin Farnese, Ricardo Figueira, Alexandre Moufarek,
386 Austin Huang, Patrik Zochbauer, Ben Ingram, Tongzhou Chen, Zelin Wu, Adrià Puigdomènec,
387 Leland Rechis, Da Yu, Sri Gayatri Sundara Padmanabhan, Rui Zhu, Chu ling Ko, Andrea Banino,
388 Samira Daruki, Aarush Selvan, Dhruva Bhaswar, Daniel Hernandez Diaz, Chen Su, Salvatore
389 Scellato, Jennifer Brennan, Woohyun Han, Grace Chung, Priyanka Agrawal, Urvashi Khandelwal,
390 Khe Chai Sim, Morgane Lustman, Sam Ritter, Kelvin Guu, Jiawei Xia, Prateek Jain, Emma Wang,
391 Tyrone Hill, Mirko Rossini, Marija Kostelac, Tautvydas Misiunas, Amit Sabne, Kyuyeun Kim,
392 Ahmet Iscen, Congchao Wang, José Leal, Ashwin Sreevatsa, Utku Evci, Manfred Warmuth, Saket
393 Joshi, Daniel Suo, James Lottes, Garrett Honke, Brendan Jou, Stefani Karp, Jieru Hu, Himanshu
394 Sahni, Adrien Ali Taïga, William Kong, Samrat Ghosh, Renshen Wang, Jay Pavagadhi, Natalie
395 Axelsson, Nikolai Grigorev, Patrick Siegler, Rebecca Lin, Guohui Wang, Emilio Parisotto, Sharath
396 Maddineni, Krishan Subudhi, Eyal Ben-David, Elena Pochevina, Orgad Keller, Thi Avrahami,
397 Zhe Yuan, Pulkit Mehta, Jialu Liu, Sherry Yang, Wendy Kan, Katherine Lee, Tom Funkhouser,
398 Derek Cheng, Hongzhi Shi, Archit Sharma, Joe Kelley, Matan Eyal, Yury Malkov, Corentin Tallec,
399 Yuval Bahat, Shen Yan, Xintian, Wu, David Lindner, Chengda Wu, Avi Caciularu, Xiyang Luo,
400 Rodolphe Jenatton, Tim Zaman, Yingying Bi, Ilya Kornakov, Ganesh Mallya, Daisuke Ikeda, Itay
401 Karo, Anima Singh, Colin Evans, Praneeth Netrapalli, Vincent Nallatamby, Isaac Tian, Yannis
402 Assael, Vikas Raunak, Victor Carbune, Ioana Bica, Lior Madmoni, Dee Cattle, Snchit Grover,
403 Krishna Somandepalli, Sid Lall, Amelio Vázquez-Reina, Riccardo Patana, Jiaqi Mu, Pranav Talluri,
404 Maggie Tran, Rajeev Aggarwal, RJ Skerry-Ryan, Jun Xu, Mike Burrows, Xiaoyue Pan, Edouard
405 Yvinec, Di Lu, Zhiying Zhang, Duc Dung Nguyen, Hairong Mu, Gabriel Barcik, Helen Ran,
406 Lauren Beltrone, Krzysztof Choromanski, Dia Kharrat, Samuel Albanie, Sean Purser-haskell,
407 David Bieber, Carrie Zhang, Jing Wang, Tom Hudson, Zhiyuan Zhang, Han Fu, Johannes Mauerer,
408 Mohammad Hossein Bateni, AJ Maschinot, Bing Wang, Muye Zhu, Arjun Pillai, Tobias Weyand,
409 Shuang Liu, Oscar Akerlund, Fred Bertsch, Vittal Premachandran, Alicia Jin, Vincent Roulet,
410 Peter de Boursac, Shubham Mittal, Ndaba Ndebele, Georgi Karadzhov, Sahra Ghalebikesabi,
411 Ricky Liang, Allen Wu, Yale Cong, Nimesh Ghelani, Sumeet Singh, Bahar Fatemi, Warren, Chen,
412 Charles Kwong, Alexey Kolganov, Steve Li, Richard Song, Chenkai Kuang, Sobhan Miryoosefi,
413 Dale Webster, James Wendt, Arkadiusz Socala, Guolong Su, Artur Mendonça, Abhinav Gupta,
414 Xiaowei Li, Tomy Tsai, Qiong, Hu, Kai Kang, Angie Chen, Sertan Girgin, Yongqin Xian, Andrew
415 Lee, Nolan Ramsden, Leslie Baker, Madeleine Clare Elish, Varvara Krayanova, Rishabh Joshi,
416 Jiri Simska, Yao-Yuan Yang, Piotr Ambroszczyk, Dipankar Ghosh, Arjun Kar, Yuan Shangguan,
417 Yumeya Yamamori, Yaroslav Akulov, Andy Brock, Haotian Tang, Siddharth Vashishtha, Rich
418 Munoz, Andreas Steiner, Kalyan Andra, Daniel Eppens, Qixuan Feng, Hayato Kobayashi, Sasha
419 Goldshtain, Mona El Mahdy, Xin Wang, Jilei, Wang, Richard Killam, Tom Kwiatkowski, Kavya
420 Kopparapu, Serena Zhan, Chao Jia, Alexei Bendebury, Sheryl Luo, Adrià Recasens, Timothy

421 Knight, Jing Chen, Mohak Patel, YaGuang Li, Ben Withbroe, Dean Weesner, Kush Bhatia, Jie
422 Ren, Danielle Eisenbud, Ebrahim Songhori, Yanhua Sun, Travis Choma, Tasos Kementsietsidis,
423 Lucas Manning, Brian Roark, Wael Farhan, Jie Feng, Susheel Tatineni, James Cobon-Kerr, Yunjie
424 Li, Lisa Anne Hendricks, Isaac Noble, Chris Breaux, Nate Kushman, Liqian Peng, Fuzhao Xue,
425 Taylor Tobin, Jamie Rogers, Josh Lipschultz, Chris Alberti, Alexey Vlaskin, Mostafa Dehghani,
426 Roshan Sharma, Tris Warkentin, Chen-Yu Lee, Benigno Urias, Da-Cheng Juan, Angad Chandorkar,
427 Hila Sheftel, Ruibo Liu, Elnaz Davoodi, Borja De Balle Pigem, Kedar Dhamdhere, David Ross,
428 Jonathan Hoech, Mahdis Mahdieh, Li Liu, Qiuja Li, Liam McCafferty, Chenxi Liu, Markus
429 Mircea, Yunting Song, Omkar Savant, Alaa Saade, Colin Cherry, Vincent Hellendoorn, Siddharth
430 Goyal, Paul Pucciarelli, David Vilar Torres, Zohar Yahav, Hyo Lee, Lars Lowe Sjoesund, Christo
431 Kirov, Bo Chang, Deepanway Ghoshal, Lu Li, Gilles Baechler, Sébastien Pereira, Tara Sainath,
432 Anudhyan Boral, Dominik Grewe, Afief Halumi, Nguyen Minh Phu, Tianxiao Shen, Marco Tulio
433 Ribeiro, Dhriti Varma, Alex Kaskasoli, Vlad Feinberg, Navneet Potti, Jarrod Kahn, Matheus
434 Wisniewski, Shakir Mohamed, Arnar Mar Hrafnelsson, Bobak Shahriari, Jean-Baptiste Lespiau,
435 Lisa Patel, Legg Yeung, Tom Paine, Lantao Mei, Alex Ramirez, Rakesh Shivanna, Li Zhong, Josh
436 Woodward, Guilherme Tubone, Samira Khan, Heng Chen, Elizabeth Nielsen, Catalin Ionescu,
437 Utsav Prabhu, Mingcen Gao, Qingze Wang, Sean Augenstein, Neesha Subramaniam, Jason Chang,
438 Fotis Iliopoulos, Jiaming Luo, Myriam Khan, Weicheng Kuo, Denis Teplyashin, Florence Perot,
439 Logan Kilpatrick, Amir Globerson, Hongkun Yu, Anfal Siddiqui, Nick Sukhanov, Arun Kandoor,
440 Umang Gupta, Marco Andreetto, Moran Ambar, Donnie Kim, Paweł Wesołowski, Sarah Perrin,
441 Ben Limonchik, Wei Fan, Jim Stephan, Ian Stewart-Binks, Ryan Kappedal, Tong He, Sarah Cogan,
442 Romina Datta, Tong Zhou, Jiayu Ye, Leandro Kieliger, Ana Ramalho, Kyle Kastner, Fabian
443 Mentzer, Wei-Jen Ko, Arun Suggala, Tianhao Zhou, Shiraz Butt, Hana Strejček, Lior Belenki,
444 Subhashini Venugopalan, Mingyang Ling, Evgenii Eltyshev, Yunxiao Deng, Geza Kovacs, Mukund
445 Raghavachari, Hanjun Dai, Tal Schuster, Steven Schwarcz, Richard Nguyen, Arthur Nguyen, Gavin
446 Buttimore, Shrestha Basu Mallick, Sudeep Gandhe, Seth Benjamin, Michal Jastrzebski, Le Yan,
447 Sugato Basu, Chris Apps, Isabel Edkins, James Allingham, Immanuel Odisho, Tomas Kociský,
448 Jewel Zhao, Linting Xue, Apoorv Reddy, Chrysovalantis Anastasiou, Aviel Atias, Sam Redmond,
449 Kieran Milan, Nicolas Heess, Herman Schmit, Allan Dafoe, Daniel Andor, Tynan Gangwani,
450 Anca Dragan, Sheng Zhang, Ashyana Kachra, Gang Wu, Siyang Xue, Kevin Aydin, Siqi Liu,
451 Yuxiang Zhou, Mahan Malihi, Austin Wu, Siddharth Gopal, Candice Schumann, Peter Stys,
452 Alek Wang, Mirek Olšák, Dangyi Liu, Christian Schallhart, Yiran Mao, Demetra Brady, Hao
453 Xu, Tomas Mery, Chawin Sitawarin, Siva Velusamy, Tom Cobley, Alex Zhai, Christian Walder,
454 Nitzan Katz, Ganesh Jawahar, Chinmay Kulkarni, Antoine Yang, Adam Paszke, Yinan Wang,
455 Bogdan Damoc, Zalán Borsos, Ray Smith, Jinning Li, Mansi Gupta, Andrei Kapishnikov, Sushant
456 Prakash, Florian Luisier, Rishabh Agarwal, Will Grathwohl, Kuangyuan Chen, Kehang Han,
457 Nikhil Mehta, Andrew Over, Shekoofeh Azizi, Lei Meng, Niccolò Dal Santo, Kelvin Zheng, Jane
458 Shapiro, Igor Petrovski, Jeffrey Hui, Amin Ghafouri, Jasper Snoek, James Qin, Mandy Jordan,
459 Caitlin Sikora, Jonathan Malmaud, Yuheng Kuang, Aga Świertlik, Ruoxin Sang, Chongyang Shi,
460 Leon Li, Andrew Rosenberg, Shubin Zhao, Andy Crawford, Jan-Thorsten Peter, Yun Lei, Xavier
461 Garcia, Long Le, Todd Wang, Julien Amelot, Dave Orr, Praneeth Kacham, Dana Alon, Gladys
462 Tyen, Abhinav Arora, James Lyon, Alex Kurakin, Mimi Ly, Theo Guidroz, Zhipeng Yan, Rina
463 Panigrahy, Pingmei Xu, Thais Kagohara, Yong Cheng, Eric Noland, Jinhyuk Lee, Jonathan Lee,
464 Cathy Yip, Maria Wang, Efrat Nehoran, Alexander Bykovsky, Zhihao Shan, Ankit Bhagatwala,
465 Chaochao Yan, Jie Tan, Guillermo Garrido, Dan Ethier, Nate Hurley, Grace Vesom, Xu Chen,
466 Siyuan Qiao, Abhishek Nayyar, Julian Walker, Paramjit Sandhu, Mihaela Rosca, Danny Swisher,
467 Mikhail Dektiarev, Josh Dillon, George-Cristian Muraru, Manuel Tragut, Artiom Myaskovsky,
468 David Reid, Marko Velic, Owen Xiao, Jasmine George, Mark Brand, Jing Li, Wenhao Yu, Shane
469 Gu, Xiang Deng, François-Xavier Aubet, Soheil Hassas Yeganeh, Fred Alcober, Celine Smith,
470 Trevor Cohn, Kay McKinney, Michael Tschanne, Ramesh Sampath, Gowoon Cheon, Liangchen
471 Luo, Luyang Liu, Jordi Orbay, Hui Peng, Gabriela Botea, Xiaofan Zhang, Charles Yoon, Cesar
472 Magalhaes, Paweł Stradomski, Ian Mackinnon, Steven Hemingray, Kumaran Venkatesan, Rhys
473 May, Jaeyoun Kim, Alex Druinsky, Jingchen Ye, Zheng Xu, Terry Huang, Jad Al Abdallah, Adil
474 Dostmohamed, Rachana Fellinger, Tsendsuren Munkhdalai, Akanksha Maurya, Peter Garst, Yin
475 Zhang, Maxim Krikun, Simon Bucher, Aditya Srikanth Veerubhotla, Yixin Liu, Sheng Li, Nishesh
476 Gupta, Jakub Adamek, Hanwen Chen, Bennett Orlando, Aleksandr Zaks, Joost van Amersfoort,
477 Josh Camp, Hui Wan, HyunJeong Choe, Zhichun Wu, Kate Olszewska, Weiren Yu, Archita Vadali,
478 Martin Scholz, Daniel De Freitas, Jason Lin, Amy Hua, Xin Liu, Frank Ding, Yichao Zhou, Boone
479 Severson, Katerina Tsihlas, Samuel Yang, Tammo Spalink, Varun Yerram, Helena Pankov, Rory

480 Blevins, Ben Vargas, Sarthak Jauhari, Matt Miecnikowski, Ming Zhang, Sandeep Kumar, Clement
481 Farabet, Charline Le Lan, Sebastian Flennerhag, Yonatan Bitton, Ada Ma, Arthur Bražinskas,
482 Eli Collins, Niharika Ahuja, Sneha Kudugunta, Anna Bortsova, Minh Giang, Wanzheng Zhu,
483 Ed Chi, Scott Lundberg, Alexey Stern, Subha Puttagunta, Jing Xiong, Xiao Wu, Yash Pande,
484 Amit Jhinal, Daniel Murphy, Jon Clark, Marc Brockschmidt, Maxine Deines, Kevin R. McKee,
485 Dan Bahir, Jiajun Shen, Minh Truong, Daniel McDuff, Andrea Gesmundo, Edouard Rosseel,
486 Bowen Liang, Ken Caluwaerts, Jessica Hamrick, Joseph Kready, Mary Cassin, Rishikesh Ingale,
487 Li Lao, Scott Pollom, Yifan Ding, Wei He, Lizzeth Bellot, Joana Iljazi, Ramya Sree Boppana,
488 Shan Han, Tara Thompson, Amr Khalifa, Anna Bulanova, Blagoj Mitrevski, Bo Pang, Emma
489 Cooney, Tian Shi, Rey Coaguila, Tamar Yakar, Marc'aurelio Ranzato, Nikola Momchev, Chris
490 Rawles, Zachary Charles, Young Maeng, Yuan Zhang, Rishabh Bansal, Xiaokai Zhao, Brian
491 Albert, Yuan Yuan, Sudheendra Vijayanarasimhan, Roy Hirsch, Vinay Ramasesh, Kiran Vodrahalli,
492 Xingyu Wang, Arushi Gupta, DJ Strouse, Jianmo Ni, Roma Patel, Gabe Taubman, Zhouyuan
493 Huo, Dero Gharibian, Marianne Monteiro, Hoi Lam, Shobha Vasudevan, Aditi Chaudhary, Isabela
494 Albuquerque, Kilol Gupta, Sebastian Riedel, Chaitra Hegde, Avraham Ruderman, András György,
495 Marcus Wainwright, Ashwin Chaugule, Burcu Karagol Ayan, Tomer Levinboim, Sam Shleifer,
496 Yogesh Kalley, Vahab Mirrokni, Abhishek Rao, Prabakar Radhakrishnan, Jay Hartford, Jialin
497 Wu, Zhenhai Zhu, Francesco Bertolini, Hao Xiong, Nicolas Serrano, Hamish Tomlinson, Myle
498 Ott, Yifan Chang, Mark Graham, Jian Li, Marco Liang, Xiangzhu Long, Sebastian Borgeaud,
499 Yanif Ahmad, Alex Grills, Diana Mincu, Martin Izzard, Yuan Liu, Jinyu Xie, Louis O'Bryan,
500 Sameera Ponda, Simon Tong, Michelle Liu, Dan Malkin, Khalid Salama, Yuankai Chen, Rohan
501 Anil, Anand Rao, Rigel Swavely, Misha Bilenko, Nina Anderson, Tat Tan, Jing Xie, Xing Wu,
502 Lijun Yu, Oriol Vinyals, Andrey Ryabtsev, Rumen Dangovski, Kate Baumli, Daniel Keysers,
503 Christian Wright, Zoe Ashwood, Betty Chan, Artem Shtefan, Yaohui Guo, Ankur Bapna, Radu
504 Soricut, Steven Pecht, Sabela Ramos, Rui Wang, Jiahao Cai, Trieu Trinh, Paul Barham, Linda
505 Friso, Eli Stickgold, Xiangzhuo Ding, Siamak Shakeri, Diego Ardila, Eleftheria Briakou, Phil
506 Culliton, Adam Raveret, Jingyu Cui, David Saxton, Subhrajit Roy, Javad Azizi, Pengcheng Yin,
507 Lucia Loher, Andrew Bunner, Min Choi, Faruk Ahmed, Eric Li, Yin Li, Shengyang Dai, Michael
508 Elabd, Sriram Ganapathy, Shivani Agrawal, Yiqing Hua, Paige Kunkle, Sujeevan Rajayogam, Arun
509 Ahuja, Arthur Conmy, Alex Vasiloff, Parker Beak, Christopher Yew, Jayaram Mudigonda, Bartek
510 Wydrowski, Jon Blanton, Zhengdong Wang, Yann Dauphin, Zhuo Xu, Martin Polacek, Xi Chen,
511 Hexiang Hu, Pauline Sho, Markus Kunesch, Mehdi Hafezi Manshadi, Eliza Rutherford, Bo Li,
512 Sissie Hsiao, Iain Barr, Alex Tudor, Matija Kecman, Arsha Nagrani, Vladimir Pchelin, Martin
513 Sundermeyer, Aishwarya P S, Abhijit Karmarkar, Yi Gao, Grishma Chole, Olivier Bachem, Isabel
514 Gao, Arturo BC, Matt Dibb, Mauro Verzetti, Felix Hernandez-Campos, Yana Lunts, Matthew
515 Johnson, Julia Di Trapani, Raphael Koster, Idan Brusilovsky, Binbin Xiong, Megha Mohabey, Han
516 Ke, Joe Zou, Tea Sabolić, Víctor Campos, John Palowitch, Alex Morris, Linhai Qiu, Pranavaraj
517 Ponnuramu, Fangtao Li, Vivek Sharma, Kiranbir Sodhia, Kaan Tekelioglu, Aleksandr Chuklin,
518 Madhavi Yenugula, Erika Gemzer, Theofilos Strinopoulos, Sam El-Husseini, Huiyu Wang, Yan
519 Zhong, Edouard Leurent, Paul Natsev, Weijun Wang, Dre Mahaarachchi, Tao Zhu, Songyou Peng,
520 Sami Alabed, Cheng-Chun Lee, Anthony Brohan, Arthur Szlam, GS Oh, Anton Kovsharov, Jenny
521 Lee, Renee Wong, Megan Barnes, Gregory Thornton, Felix Gimeno, Omer Levy, Martin Sevenich,
522 Melvin Johnson, Jonathan Mallinson, Robert Dadashi, Ziyue Wang, Qingchun Ren, Preethi Lahoti,
523 Arka Dhar, Josh Feldman, Dan Zheng, Thatcher Ulrich, Liviu Panait, Michiel Blokzijl, Cip
524 Baetu, Josip Matak, Jitendra Harlalka, Maulik Shah, Tal Marian, Daniel von Dincklage, Cosmo
525 Du, Ruy Ley-Wild, Bethanie Brownfield, Max Schumacher, Yury Stuken, Shadi Noghabi, Sonal
526 Gupta, Xiaoqi Ren, Eric Malmi, Felix Weissenberger, Blanca Huergo, Maria Bauza, Thomas
527 Lampe, Arthur Douillard, Mojtaba Seyedhosseini, Roy Frostig, Zoubin Ghahramani, Kelvin
528 Nguyen, Kashyap Krishnakumar, Chengxi Ye, Rahul Gupta, Alireza Nazari, Robert Geirhos, Pete
529 Shaw, Ahmed Eleryan, Dima Damen, Jennimaria Palomaki, Ted Xiao, Qiyin Wu, Quan Yuan,
530 Phoenix Meadowlark, Matthew Bilotti, Raymond Lin, Mukund Sridhar, Yannick Schroecker,
531 Da-Woon Chung, Jincheng Luo, Trevor Strohman, Tianlin Liu, Anne Zheng, Jesse Emond, Wei
532 Wang, Andrew Lampinen, Toshiyuki Fukuzawa, Folawiyo Campbell-Ajala, Monica Roy, James
533 Lee-Thorp, Lily Wang, Iftekhar Naim, Tony, Nguy ên, Guy Bensky, Aditya Gupta, Dominika
534 Rogozińska, Justin Fu, Thanumalayan Sankaranarayana Pillai, Petar Veličković, Shahar Drath,
535 Philipp Neubeck, Vaibhav Tulsyan, Arseniy Klimovskiy, Don Metzler, Sage Stevens, Angel
536 Yeh, Junwei Yuan, Tianhe Yu, Kelvin Zhang, Alec Go, Vincent Tsang, Ying Xu, Andy Wan,
537 Isaac Galatzer-Levy, Sam Sobell, Abodunrinwa Toki, Elizabeth Salesky, Wenlei Zhou, Diego
538 Antognini, Sholto Douglas, Shimu Wu, Adam Lelkes, Frank Kim, Paul Cavallaro, Ana Salazar,

539 Yuchi Liu, James Besley, Tiziana Refice, Yiling Jia, Zhang Li, Michal Sokolik, Arvind Kannan,
540 Jon Simon, Jo Chick, Avia Aharon, Meet Gandhi, Mayank Daswani, Keyvan Amiri, Vighnesh
541 Birodkar, Abe Ittycheriah, Peter Grabowski, Oscar Chang, Charles Sutton, Zhixin, Lai, Umesh
542 Telang, Susie Sargsyan, Tao Jiang, Raphael Hoffmann, Nicole Brichtova, Matteo Hessel, Jonathan
543 Halcrow, Sammy Jerome, Geoff Brown, Alex Tomala, Elena Buchatskaya, Dian Yu, Sachit
544 Menon, Pol Moreno, Yuguo Liao, Vicky Zayats, Luming Tang, SQ Mah, Ashish Shenoy, Alex
545 Siegman, Majid Hadian, Okwan Kwon, Tao Tu, Nima Khajehnouri, Ryan Foley, Parisa Haghani,
546 Zhongru Wu, Vaishakh Keshava, Khyatti Gupta, Tony Bruguier, Rui Yao, Danny Karmon, Luisa
547 Zintgraf, Zhicheng Wang, Enrique Piqueras, Junehyuk Jung, Jenny Brennan, Diego Machado,
548 Marissa Giustina, MH Tessler, Kamyu Lee, Qiao Zhang, Joss Moore, Kaspar Daugaard, Alexander
549 Frömmgen, Jennifer Beattie, Fred Zhang, Daniel Kasenberg, Ty Geri, Danfeng Qin, Gaurav Singh
550 Tomar, Tom Ouyang, Tianli Yu, Luowei Zhou, Rajiv Mathews, Andy Davis, Yaoyiran Li, Jai
551 Gupta, Damion Yates, Linda Deng, Elizabeth Kemp, Ga-Young Joung, Sergei Vassilvitskii, Mandy
552 Guo, Pallavi LV, Dave Dopson, Sami Lachgar, Lara McConaughey, Himadri Choudhury, Dragos
553 Dena, Aaron Cohen, Joshua Ainslie, Sergey Levi, Parthasarathy Gopavarapu, Polina Zablotskaia,
554 Hugo Vallet, Sanaz Bahargam, Xiaodan Tang, Nenad Tomasev, Ethan Dyer, Daniel Balle, Hongrae
555 Lee, William Bono, Jorge Gonzalez Mendez, Vadim Zubov, Shentao Yang, Ivor Rendulic, Yanyan
556 Zheng, Andrew Hogue, Golan Pundak, Ralph Leith, Avishkar Bhoopchand, Michael Han, Mislav
557 Žanić, Tom Schaul, Manolis Delakis, Tejas Iyer, Guanyu Wang, Harman Singh, Abdelrahman
558 Abdelhamed, Tara Thomas, Siddhartha Brahma, Hilal Dib, Naveen Kumar, Wenxuan Zhou, Liang
559 Bai, Pushkar Mishra, Jiao Sun, Valentin Anklin, Roykrong Sukerd, Lauren Agubuzu, Anton
560 Briukhov, Anmol Gulati, Maximilian Sieb, Fabio Pardo, Sara Nasso, Junquan Chen, Kexin Zhu,
561 Tiberiu Sosea, Alex Goldin, Keith Rush, Spurthi Amba Hombaiah, Andreas Noever, Allan Zhou,
562 Sam Haves, Mary Phuong, Jake Ades, Yi ting Chen, Lin Yang, Joseph Pagadora, Stan Bileschi,
563 Victor Cotruta, Rachel Saputro, Arijit Pramanik, Sean Ammirati, Dan Garrette, Kevin Villela, Tim
564 Blyth, Canfer Akbulut, Neha Jha, Alban Rustemi, Arissa Wongpanich, Chirag Nagpal, Yonghui
565 Wu, Morgane Rivière, Sergey Kishchenko, Pranesh Srinivasan, Alice Chen, Animesh Sinha, Trang
566 Pham, Bill Jia, Tom Hennigan, Anton Bakalov, Nithya Attaluri, Drew Garmon, Daniel Rodriguez,
567 Dawid Wegner, Wenhao Jia, Evan Senter, Noah Fiedel, Denis Petek, Yuchuan Liu, Cassidy Hardin,
568 Harshal Tushar Lehri, Joao Carreira, Sara Smoot, Marcel Prasetya, Nami Akazawa, Anca Stefanou,
569 Chia-Hua Ho, Anelia Angelova, Kate Lin, Min Kim, Charles Chen, Marcin Sieniek, Alice Li,
570 Tongfei Guo, Sorin Baltateanu, Pouya Tafti, Michael Wunder, Nadav Olmert, Divyansh Shukla,
571 Jingwei Shen, Neel Kovelamudi, Balaji Venkatraman, Seth Neel, Romal Thoppilan, Jerome Connor,
572 Frederik Benzing, Axel Stjerngren, Golnaz Ghiasi, Alex Polozov, Joshua Howland, Theophane
573 Weber, Justin Chiu, Ganesh Poomal Girirajan, Andreas Terzis, Pidong Wang, Fangda Li, Yoav Ben
574 Shalom, Dinesh Tewari, Matthew Denton, Roee Aharoni, Norbert Kalb, Heri Zhao, Junlin Zhang,
575 Angelos Filos, Matthew Rahtz, Lalit Jain, Connie Fan, Vitor Rodrigues, Ruth Wang, Richard
576 Shin, Jacob Austin, Roman Ring, Mariella Sanchez-Vargas, Mehadi Hassen, Ido Kessler, Uri Alon,
577 Gufeng Zhang, Wenhua Chen, Yenai Ma, Xiance Si, Le Hou, Azalia Mirhoseini, Marc Wilson,
578 Geoff Bacon, Becca Roelofs, Lei Shu, Gautam Vasudevan, Jonas Adler, Artur Dwornik, Tayfun
579 Terzi, Matt Lawlor, Harry Askham, Mike Bernico, Xuanyi Dong, Chris Hidey, Kevin Kilgour,
580 Gaël Liu, Surya Bhupatiraju, Luke Leonhard, Siqi Zuo, Partha Talukdar, Qing Wei, Aliaksei
581 Severyn, Vít Listík, Jong Lee, Aditya Tripathi, SK Park, Yossi Matias, Hao Liu, Alex Ruiz, Rajesh
582 Jayaram, Jackson Tolins, Pierre Marcenac, Yiming Wang, Bryan Seybold, Henry Prior, Deepak
583 Sharma, Jack Weber, Mikhail Sirotenko, Yunhsuan Sung, Dayou Du, Ellie Pavlick, Stefan Zinke,
584 Markus Freitag, Max Dylla, Montse Gonzalez Arenas, Natan Potikha, Omer Goldman, Connie
585 Tao, Rachita Chhaparia, Maria Voitovich, Pawan Dogra, Andrija Ražnatović, Zak Tsai, Chong
586 You, Oleaser Johnson, George Tucker, Chenjie Gu, Jae Yoo, Maryam Majzoubi, Valentin Gabeur,
587 Bahram Raad, Rocky Rhodes, Kashyap Kolipaka, Heidi Howard, Geta Sampemane, Benny Li,
588 Chulayuth Asawaroengchai, Duy Nguyen, Chiyuan Zhang, Timothee Cour, Xinxin Yu, Zhao Fu,
589 Joe Jiang, Po-Sen Huang, Gabriela Surita, Iñaki Iturrate, Yael Karov, Michael Collins, Martin
590 Baeuml, Fabian Fuchs, Shilpa Shetty, Swaroop Ramaswamy, Sayna Ebrahimi, Qiuchen Guo,
591 Jeremy Shar, Gabe Barth-Maron, Sravanti Addepalli, Bryan Richter, Chin-Yi Cheng, Eugénie
592 Rives, Fei Zheng, Johannes Griesser, Nishanth Dikkala, Yoel Zeldes, Ilkin Safarli, Dipanjan Das,
593 Himanshu Srivastava, Sadh MNM Khan, Xin Li, Aditya Pandey, Larisa Markeeva, Dan Belov, Qiqi
594 Yan, Mikołaj Rybiński, Tao Chen, Megha Nawhal, Michael Quinn, Vineetha Govindaraj, Sarah
595 York, Reed Roberts, Roopal Garg, Namrata Godbole, Jake Abernethy, Anil Das, Lam Nguyen
596 Thiet, Jonathan Tompson, John Nham, Neera Vats, Ben Caine, Wesley Helmholz, Francesco
597 Pongetti, Yeongil Ko, James An, Clara Huiyi Hu, Yu-Cheng Ling, Julia Pawar, Robert Leland,

598 Keisuke Kinoshita, Waleed Khawaja, Marco Selvi, Eugene Ie, Danila Sinopalnikov, Lev Proleev,
599 Nilesh Tripuraneni, Michele Bevilacqua, Seungji Lee, Clayton Sanford, Dan Suh, Dustin Tran,
600 Jeff Dean, Simon Baumgartner, Jens Heitkaemper, Sagar Gubbi, Kristina Toutanova, Yichong Xu,
601 Chandu Thekkath, Keran Rong, Palak Jain, Annie Xie, Yan Virin, Yang Li, Lubo Litchev, Richard
602 Powell, Tarun Bharti, Adam Kraft, Nan Hua, Marissa Ikonomidis, Ayal Hitron, Sanjiv Kumar,
603 Loic Matthey, Sophie Bridgers, Lauren Lax, Ishaan Malhi, Ondrej Skopek, Ashish Gupta, Jiawei
604 Cao, Mitchelle Rasquinha, Siim Pöder, Wojciech Stokowiec, Nicholas Roth, Guowang Li, Michaël
605 Sander, Joshua Kessinger, Vihan Jain, Edward Loper, Wonpyo Park, Michal Yarom, Lijun Cheng,
606 Guru Guruganesh, Kanishka Rao, Yan Li, Catarina Barros, Mikhail Sushkov, Chun-Sung Ferng,
607 Rohin Shah, Ophir Aharoni, Ravin Kumar, Tim McConnell, Peiran Li, Chen Wang, Fernando
608 Pereira, Craig Swanson, Fayaz Jamil, Yan Xiong, Anitha Vijayakumar, Prakash Shroff, Kedar
609 Soparkar, Jindong Gu, Livio Baldini Soares, Eric Wang, Kushal Majmundar, Aurora Wei, Kai
610 Bailey, Nora Kassner, Chizu Kawamoto, Goran Žužić, Victor Gomes, Abhirut Gupta, Michael
611 Guzman, Ishita Dasgupta, Xinyi Bai, Zhufeng Pan, Francesco Piccinno, Hadas Natalie Vogel,
612 Octavio Ponce, Adrian Hutter, Paul Chang, Pan-Pan Jiang, Ionel Gog, Vlad Ionescu, James
613 Manyika, Fabian Pedregosa, Harry Ragan, Zach Behrman, Ryan Mullins, Coline Devin, Aroonalok
614 Pyne, Swapnil Gawde, Martin Chadwick, Yiming Gu, Sasan Tavakkol, Andy Twigg, Naman
615 Goyal, Ndidi Elue, Anna Goldie, Srinivasan Venkatachary, Hongliang Fei, Ziqiang Feng, Marvin
616 Ritter, Isabel Leal, Sudeep Dasari, Pei Sun, Alif Raditya Rochman, Brendan O'Donoghue, Yuchen
617 Liu, Jim Sproch, Kai Chen, Natalie Clay, Slav Petrov, Sailesh Sidhwani, Ioana Mihailescu, Alex
618 Panagopoulos, AJ Piergiovanni, Yunfei Bai, George Powell, Deep Karkhanis, Trevor Yacovone,
619 Petr Mitrichev, Joe Kovac, Dave Uthus, Amir Yazdanbakhsh, David Amos, Steven Zheng, Bing
620 Zhang, Jin Miao, Bhuvana Ramabhadran, Soroush Radpour, Shantanu Thakoor, Josh Newlan, Oran
621 Lang, Orion Jankowski, Shikhar Bharadwaj, Jean-Michel Sarr, Shereen Ashraf, Sneha Mondal, Jun
622 Yan, Ankit Singh Rawat, Sarmishta Velury, Greg Kochanski, Tom Eccles, Franz Och, Abhanshu
623 Sharma, Ethan Mahintorabi, Alex Gurney, Carrie Muir, Vered Cohen, Saksham Thakur, Adam
624 Bloniarz, Asier Mujika, Alexander Pritzel, Paul Caron, Altaf Rahman, Fiona Lang, Yasumasa Onoe,
625 Petar Sirkovic, Jay Hoover, Ying Jian, Pablo Duque, Arun Narayanan, David Soergel, Alex Haig,
626 Loren Maggiore, Shyamal Buch, Josef Dean, Ilya Figotin, Igor Karpov, Shaleen Gupta, Denny
627 Zhou, Muhan Huang, Ashwin Vaswani, Christopher Semturs, Kaushik Shivakumar, Yu Watanabe,
628 Vinodh Kumar Rajendran, Eva Lu, Yanhan Hou, Wenting Ye, Shikhar Vashishth, Nana Nti, Vytenis
629 Sakenas, Darren Ni, Doug DeCarlo, Michael Bendersky, Sumit Bagri, Nacho Cano, Elijah Peake,
630 Simon Tokumine, Varun Godbole, Carlos Guía, Tanya Lando, Vittorio Selo, Seher Ellis, Danny
631 Tarlow, Daniel Gillick, Alessandro Epasto, Siddhartha Reddy Jonnalagadda, Meng Wei, Meiyang
632 Xie, Ankur Taly, Michela Paganini, Mukund Sundararajan, Daniel Toyama, Ting Yu, Dessie
633 Petrova, Aneesh Pappu, Rohan Agrawal, Senaka Butphitiya, Justin Frye, Thomas Buschmann,
634 Remi Crocker, Marco Tagliasacchi, Mengchao Wang, Da Huang, Sagi Perel, Brian Wieder, Hideto
635 Kazawa, Weiyue Wang, Jeremy Cole, Himanshu Gupta, Ben Golan, Seojin Bang, Nitish Kulkarni,
636 Ken Franko, Casper Liu, Doug Reid, Sid Dalmia, Jay Whang, Kevin Cen, Prasha Sundaram, Johan
637 Ferret, Berivan Isik, Lucian Ionita, Guan Sun, Anna Shekhawat, Muqthar Mohammad, Philip
638 Pham, Ronny Huang, Karthik Raman, Xingyi Zhou, Ross McIlroy, Austin Myers, Sheng Peng,
639 Jacob Scott, Paul Covington, Sofia Erell, Pratik Joshi, João Gabriel Oliveira, Natasha Noy, Tajwar
640 Nasir, Jake Walker, Vera Axelrod, Tim Dozat, Pu Han, Chun-Te Chu, Eugene Weinstein, Anand
641 Shukla, Shreyas Chandrakaladharan, Petra Poklukar, Bonnie Li, Ye Jin, Prem Erubetidine, Steven
642 Hansen, Avigail Dabush, Alon Jacovi, Samrat Phatale, Chen Zhu, Steven Baker, Mo Shomrat, Yang
643 Xiao, Jean Pouget-Abadie, Mingyang Zhang, Fanny Wei, Yang Song, Helen King, Yiling Huang,
644 Yun Zhu, Ruoxi Sun, Juliana Vicente Franco, Chu-Cheng Lin, Sho Arora, Hui, Li, Vivian Xia,
645 Luke Vilnis, Mariano Schain, Kaiz Alarakyia, Laurel Prince, Aaron Phillips, Caleb Habtegebriel,
646 Luyao Xu, Huan Gui, Santiago Ontanon, Lora Aroyo, Karan Gill, Peggy Lu, Yash Katariya,
647 Dhruv Madeka, Shankar Krishnan, Shubha Srinivas Raghvendra, James Freedman, Yi Tay, Gaurav
648 Menghani, Peter Choy, Nishita Shetty, Dan Abolafia, Doron Kukliansky, Edward Chou, Jared
649 Lichtarge, Ken Burke, Ben Coleman, Dee Guo, Larry Jin, Indro Bhattacharya, Victoria Langston,
650 Yiming Li, Suyog Kotecha, Alex Yakubovich, Xinyun Chen, Petre Petrov, Tolly Powell, Yanzhang
651 He, Corbin Quick, Kanav Garg, Dawsen Hwang, Yang Lu, Srinadh Bhojanapalli, Kristian Kjems,
652 Ramin Mehran, Aaron Archer, Hado van Hasselt, Ashwin Balakrishna, JK Kearns, Meiqi Guo,
653 Jason Riesa, Mikita Sazanovich, Xu Gao, Chris Sauer, Chengrun Yang, XiangHai Sheng, Thomas
654 Jimma, Wouter Van Gansbeke, Vitaly Nikolaev, Wei Wei, Katie Millican, Ruizhe Zhao, Justin
655 Snyder, Levent Bolelli, Maura O'Brien, Shawn Xu, Fei Xia, Wentao Yuan, Arvind Neelakantan,
656 David Barker, Sachin Yadav, Hannah Kirkwood, Farooq Ahmad, Joel Wee, Jordan Grimstad, Boyu

657 Wang, Matthew Wiethoff, Shane Settle, Miaosen Wang, Charles Blundell, Jingjing Chen, Chris
658 Duvarney, Grace Hu, Olaf Ronneberger, Alex Lee, Yuanzhen Li, Abhishek Chakladar, Alena
659 Butryna, Georgios Evangelopoulos, Guillaume Desjardins, Jonni Kanerva, Henry Wang, Averi
660 Nowak, Nick Li, Alyssa Loo, Art Khurshudov, Laurent El Shafey, Nagabhushan Baddi, Karel Lenc,
661 Yasaman Razeghi, Tom Lieber, Amer Sinha, Xiao Ma, Yao Su, James Huang, Asahi Ushio, Hanna
662 Klimczak-Plucińska, Kareem Mohamed, JD Chen, Simon Osindero, Stav Ginzburg, Lampros
663 Lamprou, Vasilisa Bashlovkina, Duc-Hieu Tran, Ali Khodaei, Ankit Anand, Yixian Di, Ramy
664 Eskander, Manish Reddy Vuyyuru, Jasmine Liu, Aishwarya Kamath, Roman Goldenberg, Mathias
665 Bellaiche, Juliette Pluto, Bill Rosgen, Hassan Mansoor, William Wong, Suhas Ganesh, Eric Bailey,
666 Scott Baird, Dan Deutsch, Jinoo Baek, Xuhui Jia, Chansoo Lee, Abe Friesen, Nathaniel Braun, Kate
667 Lee, Amayika Panda, Steven M. Hernandez, Duncan Williams, Jianqiao Liu, Ethan Liang, Arnaud
668 Autef, Emily Pitler, Deepali Jain, Phoebe Kirk, Oskar Bunyan, Jaume Sanchez Elias, Tongxin Yin,
669 Machel Reid, Aedan Pope, Nikita Putikhin, Bidisha Samanta, Sergio Guadarrama, Dahun Kim,
670 Simon Rowe, Marcella Valentine, Geng Yan, Alex Salcianu, David Silver, Gan Song, Richa Singh,
671 Shuai Ye, Hannah DeBalsi, Majd Al Merey, Eran Ofek, Albert Webson, Shibli Mourad, Ashwin
672 Kakarla, Silvio Lattanzi, Nick Roy, Evgeny Sluzhaev, Christina Butterfield, Alessio Tonioni,
673 Nathan Waters, Sudhindra Kopalle, Jason Chase, James Cohan, Girish Ramchandra Rao, Robert
674 Berry, Michael Voznesensky, Shuguang Hu, Kristen Chiafullo, Sharat Chikkerur, George Scrivener,
675 Ivy Zheng, Jeremy Wiesner, Wolfgang Macherey, Timothy Lillicrap, Fei Liu, Brian Walker, David
676 Welling, Elinor Davies, Yangsibo Huang, Lijie Ren, Nir Shabat, Alessandro Agostini, Mariko
677 Iinuma, Dustin Zelle, Rohit Sathyaranayana, Andrea D'olimpio, Morgan Redshaw, Matt Ginsberg,
678 Ashwin Murthy, Mark Geller, Tatiana Matejovicova, Ayan Chakrabarti, Ryan Julian, Christine
679 Chan, Qiong Hu, Daniel Jarrett, Manu Agarwal, Jeshwanth Challagundla, Tao Li, Sandeep Tata,
680 Wen Ding, Maya Meng, Zhuyun Dai, Giulia Vezzani, Shefali Garg, Jannis Bulian, Mary Jasarevic,
681 Honglong Cai, Harish Rajamani, Adam Santoro, Florian Hartmann, Chen Liang, Bartek Perz,
682 Apoorv Jindal, Fan Bu, Sungyong Seo, Ryan Poplin, Adrian Goedeckemeyer, Badih Ghazi, Nikhil
683 Khadke, Leon Liu, Kevin Mather, Mingda Zhang, Ali Shah, Alex Chen, Jinliang Wei, Keshav
684 Shivam, Yuan Cao, Donghyun Cho, Angelo Scorza Scarpati, Michael Moffitt, Clara Barbu, Ivan
685 Jurin, Ming-Wei Chang, Hongbin Liu, Hao Zheng, Shachi Dave, Christine Kaeser-Chen, Xiaobin
686 Yu, Alvin Abdagic, Lucas Gonzalez, Yanping Huang, Peilin Zhong, Cordelia Schmid, Bryce
687 Petri, Alex Wertheim, Jifan Zhu, Hoang Nguyen, Kaiyang Ji, Yanqi Zhou, Tao Zhou, Fangxiaoyu
688 Feng, Regev Cohen, David Rim, Shubham Milind Phal, Petko Georgiev, Ariel Brand, Yue Ma,
689 Wei Li, Somit Gupta, Chao Wang, Pavel Dubov, Jean Tarbouriech, Kingshuk Majumder, Huijian
690 Li, Norman Rink, Apurv Suman, Yang Guo, Yinghao Sun, Arun Nair, Xiaowei Xu, Mohamed
691 Elhawaty, Rodrigo Cabrera, Guangxing Han, Julian Eisenschlos, Junwen Bai, Yuqi Li, Yamini
692 Bansal, Thibault Sellam, Mina Khan, Hung Nguyen, Justin Mao-Jones, Nikos Parotsidis, Jake
693 Marcus, Cindy Fan, Roland Zimmermann, Yony Kochinski, Laura Graesser, Feryal Behbahani,
694 Alvaro Caceres, Michael Riley, Patrick Kane, Sandra Lefdal, Rob Willoughby, Paul Vicol, Lun
695 Wang, Shujian Zhang, Ashleah Gill, Yu Liang, Gautam Prasad, Soroosh Mariooryad, Mehran
696 Kazemi, Zifeng Wang, Kritika Muralidharan, Paul Voigtlaender, Jeffrey Zhao, Huanjie Zhou,
697 Nina D'Souza, Aditi Mavalankar, Séb Arnold, Nick Young, Obaid Sarvana, Chace Lee, Milad
698 Nasr, Tingting Zou, Seokhwan Kim, Lukas Haas, Kaushal Patel, Neslihan Bulut, David Parkinson,
699 Courtney Biles, Dmitry Kalashnikov, Chi Ming To, Aviral Kumar, Jessica Austin, Alex Greve,
700 Lei Zhang, Megha Goel, Yeqing Li, Sergey Yaroshenko, Max Chang, Abhishek Jindal, Geoff
701 Clark, Hagai Taitelbaum, Dale Johnson, Ofir Roval, Jeongwoo Ko, Anhad Mohananey, Christian
702 Schuler, Shenil Dodhia, Ruichao Li, Kazuki Osawa, Claire Cui, Peng Xu, Rushin Shah, Tao Huang,
703 Ela Gruzewska, Nathan Clement, Mudit Verma, Olcan Sercinoglu, Hai Qian, Viral Shah, Masa
704 Yamaguchi, Abhinit Modi, Takahiro Kosakai, Thomas Strohmann, Junhao Zeng, Beliz Gunel, Jun
705 Qian, Austin Tarango, Krzysztof Jastrzębski, Robert David, Jyn Shan, Parker Schuh, Kunal Lad,
706 Willi Gierke, Mukundan Madhavan, Xinyi Chen, Mark Kurzeja, Rebeca Santamaria-Fernandez,
707 Dawn Chen, Alexandra Cordell, Yuri Chervonyi, Frankie Garcia, Nithish Kannen, Vincent Perot,
708 Nan Ding, Shlomi Cohen-Ganor, Victor Lavrenko, Junru Wu, Georgie Evans, Cicero Nogueira dos
709 Santos, Madhavi Sewak, Ashley Brown, Andrew Hard, Joan Puigcerver, Zeyu Zheng, Yizhong
710 Liang, Evgeny Gladchenko, Reeve Ingle, Uri First, Pierre Sermanet, Charlotte Magister, Mihajlo
711 Velimirović, Sashank Reddi, Susanna Ricco, Eirikur Agustsson, Hartwig Adam, Nir Levine, David
712 Gaddy, Dan Holtmann-Rice, Xuanhui Wang, Ashutosh Sathe, Abhijit Guha Roy, Blaž Bratanič,
713 Alen Carin, Harsh Mehta, Silvano Bonacina, Nicola De Cao, Mara Finkelstein, Verena Rieser,
714 Xinyi Wu, Florent Altché, Dylan Scandinaro, Li Li, Nino Vieillard, Nikhil Sethi, Garrett Tanzer,
715 Zhi Xing, Shibo Wang, Parul Bhatia, Gui Citovsky, Thomas Anthony, Sharon Lin, Tianze Shi,

- 716 Shoshana Jakobovits, Gena Gibson, Raj Apte, Lisa Lee, Mingqing Chen, Arunkumar Byravan,
717 Petros Maniatis, Kellie Webster, Andrew Dai, Pu-Chin Chen, Jiaqi Pan, Asya Fadeeva, Zach
718 Gleicher, Thang Luong, and Niket Kumar Bhumihar. Gemini 2.5: Pushing the frontier with
719 advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025.
720 URL <https://arxiv.org/abs/2507.06261>.
- 721 Chunyuan Deng, Yilun Zhao, Yuzhao Heng, Yitong Li, Jiannan Cao, Xiangru Tang, and Arman
722 Cohan. Unveiling the spectrum of data contamination in language models: A survey from detection
723 to remediation. *arXiv preprint arXiv:2406.14644*, 2024a.
- 724 Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerstein, and Arman Cohan. Investigating data
725 contamination in modern benchmarks for large language models. In Kevin Duh, Helena Gomez,
726 and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter
727 of the Association for Computational Linguistics: Human Language Technologies (Volume 1:
728 Long Papers)*, pp. 8706–8719, Mexico City, Mexico, June 2024b. Association for Computational
729 Linguistics. doi: 10.18653/v1/2024.nacl-long.482. URL [https://aclanthology.org/2024.nacl-long.482/](https://aclanthology.org/2024.nacl-long.482).
- 730
- 731 Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld,
732 Margaret Mitchell, and Matt Gardner. Documenting large webtext corpora: A case study on the
733 colossal clean crawled corpus. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and
734 Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural
735 Language Processing*, pp. 1286–1305, Online and Punta Cana, Dominican Republic, November
736 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.98. URL
737 <https://aclanthology.org/2021.emnlp-main.98/>.
- 738 Yihong Dong, Xue Jiang, Huanyu Liu, Zhi Jin, Bin Gu, Mengfei Yang, and Ge Li. Generalization
739 or memorization: Data contamination and trustworthy evaluation for large language models.
740 In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for
741 Computational Linguistics: ACL 2024*, pp. 12039–12050, Bangkok, Thailand, August 2024.
742 Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.716. URL <https://aclanthology.org/2024.findings-acl.716/>.
- 743
- 744 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim
745 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language models
746 with mixture-of-experts. In *International conference on machine learning*, pp. 5547–5569. PMLR,
747 2022.
- 748 Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer,
749 Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do membership inference
750 attacks work on large language models? *arXiv preprint arXiv:2402.07841*, 2024.
- 751 Sunny Duan, Mikail Khona, Abhiram Iyer, Rylan Schaeffer, and Ila R Fiete. Uncovering latent
752 memories in large language models. In *The Thirteenth International Conference on Learning
753 Representations*, 2025. URL <https://openreview.net/forum?id=KSBx6FBZpE>.
- 754 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,
755 Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff,
756 Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika,
757 Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation
758 harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- 759 Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caro-
760 line Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, Olli Järvinieniemi,
761 Matthew Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla, Qiuyu Ren, Elizabeth Pratt, Lionel
762 Levine, Grant Barkley, Natalie Stewart, Bogdan Grechuk, Tetiana Grechuk, Shreepranav Varma
763 Enugandla, and Mark Wildon. Frontiermath: A benchmark for evaluating advanced mathematical
764 reasoning in ai, 2025. URL <https://arxiv.org/abs/2411.04872>.
- 765 Shahriar Golchin and Mihai Surdeanu. Data contamination quiz: A tool to detect and estimate
766 contamination in large language models. *arXiv preprint arXiv:2311.06233*, 2023.

- 767 Shahriar Golchin and Mihai Surdeanu. Time travel in LLMs: Tracing data contamination in large
768 language models. In *The Twelfth International Conference on Learning Representations*, 2024.
769 URL <https://openreview.net/forum?id=2Rwq6c3tvr>.
- 770 Ziwen Han, Meher Mankikar, Julian Michael, and Zifan Wang. Search-time data contamination,
771 2025. URL <https://arxiv.org/abs/2508.13180>.
- 772 Jamie Hayes, Marika Swanberg, Harsh Chaudhari, Itay Yona, Ilia Shumailov, Milad Nasr, Christo-
773 pher A Choquette-Choo, Katherine Lee, and A Feder Cooper. Measuring memorization in language
774 models via probabilistic extraction. In *Proceedings of the 2025 Conference of the Nations of the*
775 *Americas Chapter of the Association for Computational Linguistics: Human Language Technolo-*
776 *gies (Volume 1: Long Papers)*, pp. 9266–9291, 2025.
- 777 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
778 Steinhardt. Measuring massive multitask language understanding. In *International Conference on*
779 *Learning Representations*, 2021a.
- 780 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
781 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021b.
782 URL <https://arxiv.org/abs/2103.03874>.
- 783 Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson
784 Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, Scott Johnston, Ben Mann, Chris Olah,
785 Catherine Olsson, Dario Amodei, Nicholas Joseph, Jared Kaplan, and Sam McCandlish. Scaling
786 laws and interpretability of learning from repeated data, 2022. URL <https://arxiv.org/abs/2205.10487>.
- 787 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
788 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas
789 Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Au-
790 relia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and Lau-
791 rent Sifre. An empirical analysis of compute-optimal large language model training. In
792 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances*
793 *in Neural Information Processing Systems*, volume 35, pp. 30016–30030. Curran Associates,
794 Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faf6f588870935f114ebe04a3e5-Paper-Conference.pdf.
- 795 Matthew Jagielski, Milad Nasr, Katherine Lee, Christopher A Choquette-Choo, Nicholas Carlini,
796 and Florian Tramer. Students parrot their teachers: Membership inference on model distillation.
797 *Advances in Neural Information Processing Systems*, 36, 2024.
- 798 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
799 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
800 evaluation of large language models for code. In *The Thirteenth International Conference on*
801 *Learning Representations*, 2025. URL <https://openreview.net/forum?id=chfJJYC3iL>.
- 802 Mingjian Jiang, Ken Ziyu Liu, and Sanmi Koyejo. A missing testbed for LLM pre-training member-
803 ship inference attacks. In *ICLR 2025 Workshop on Navigating and Addressing Data Problems for*
804 *Foundation Models*, 2025. URL <https://openreview.net/forum?id=HzHUxo6KzE>.
- 805 Minhao Jiang, Ken Ziyu Liu, Ming Zhong, Rylan Schaeffer, Siru Ouyang, Jiawei Han, and Sanmi
806 Koyejo. Investigating data contamination for pre-training language models, 2024. URL <https://arxiv.org/abs/2401.06059>.
- 807 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R
808 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth*
809 *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VTF8yNQM66>.
- 810 Nikhil Kandpal, Krishna Pillutla, Alina Oprea, Peter Kairouz, Christopher A Choquette-Choo, and
811 Zheng Xu. User inference attacks on large language models. *arXiv preprint arXiv:2310.09266*,
812 2023.

- 817 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
818 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
819 *arXiv preprint arXiv:2001.08361*, 2020.
- 820 Muhammed Yusuf Kocyigit, Eleftheria Briakou, Daniel Deutsch, Jiaming Luo, Colin Cherry, and
821 Markus Freitag. Overestimation in llm evaluation: A controlled large-scale study on data contami-
822 nation’s impact on machine translation. In *Forty-second International Conference on Machine
823 Learning*, 2025.
- 824 Zhifeng Kong, Amrita Roy Chowdhury, and Kamalika Chaudhuri. Can membership inferencing be
825 refuted? *arXiv preprint arXiv:2303.03648*, 2023.
- 826 Huihan Li, You Chen, Siyuan Wang, Yixin He, Ninareh Mehrabi, Rahul Gupta, and Xiang Ren.
827 Diagnosing memorization in chain-of-thought reasoning, one token at a time. *arXiv preprint
828 arXiv:2508.02037*, 2025.
- 829 Marvin Li, Jason Wang, Jeffrey Wang, and Seth Neel. Mope: Model perturbation-based privacy
830 attacks on language models. *arXiv preprint arXiv:2310.14369*, 2023.
- 831 Yucheng Li, Yunhao Guo, Frank Guerin, and Chenghua Lin. An open-source data contamination
832 report for large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.),
833 *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 528–541, Miami,
834 Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/
835 2024.findings-emnlp.30. URL <https://aclanthology.org/2024.findings-emnlp.30/>.
- 836 Ken Ziyu Liu, Christopher A Choquette-Choo, Matthew Jagielski, Peter Kairouz, Sanmi Koyejo,
837 Percy Liang, and Nicolas Papernot. Language models may verbatim complete text they were not
838 explicitly trained on. *arXiv preprint arXiv:2503.17514*, 2025.
- 839 Xingyu Lu, Xiaonan Li, Qinyuan Cheng, Kai Ding, Xuan-Jing Huang, and Xipeng Qiu. Scaling laws
840 for fact memorization of large language models. In *Findings of the Association for Computational
841 Linguistics: EMNLP 2024*, pp. 11263–11282, 2024.
- 842 Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. *arXiv
843 preprint arXiv:2203.08242*, 2022.
- 844 Pratyush Maini, Mohammad Yaghini, and Nicolas Papernot. Dataset inference: Ownership resolution
845 in machine learning. *arXiv preprint arXiv:2104.10706*, 2021.
- 846 Pratyush Maini, Hengrui Jia, Nicolas Papernot, and Adam Dziedzic. Llm dataset inference: Did you
847 train on my dataset? *arXiv preprint arXiv:2406.06443*, 2024.
- 848 Neal Mangaokar, Ashish Hooda, Zhuohang Li, Bradley A Malin, Kassem Fawaz, Somesh Jha,
849 Atul Prakash, and Amrita Roy Chowdhury. What really is a member? discrediting membership
850 inference via poisoning. *arXiv preprint arXiv:2506.06003*, 2025.
- 851 Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schölkopf, Mrinmaya Sachan,
852 and Taylor Berg-Kirkpatrick. Membership inference attacks against language models via neigh-
853 bourhood comparison. *arXiv preprint arXiv:2305.18462*, 2023.
- 854 Alexandre Matton, Tom Sherborne, Dennis Aumiller, Elena Tommasone, Milad Alizadeh, Jingyi He,
855 Raymond Ma, Maxime Voisin, Ellen Gilsean-McMahon, and Matthias Gallé. On leakage of code
856 generation evaluation datasets. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.),
857 *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 13215–13223, Mi-
858 ami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/
859 2024.findings-emnlp.772. URL <https://aclanthology.org/2024.findings-emnlp.772/>.
- 860 Matthieu Meeus, Shubham Jain, Marek Rei, and Yves-Alexandre de Montjoye. Inherent challenges
861 of post-hoc membership inference for large language models. *arXiv preprint arXiv:2406.17975*,
862 2024.

- 863 Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. Abstractive
864 text summarization using sequence-to-sequence RNNs and beyond. In Stefan Riezler and Yoav
865 Goldberg (eds.), *Proceedings of the 20th SIGNLL Conference on Computational Natural Language*
866 *Learning*, pp. 280–290, Berlin, Germany, August 2016. Association for Computational Linguistics.
867 doi: 10.18653/v1/K16-1028. URL <https://aclanthology.org/K16-1028/>.
- 868 Shiwen Ni, Xiangtao Kong, Chengming Li, Xiping Hu, Ruifeng Xu, Jia Zhu, and Min Yang. Training
869 on the benchmark is not all you need, 2025. URL <https://arxiv.org/abs/2409.01790>.
- 870 Fan Nie, Ken Ziyu Liu, Zihao Wang, Rui Sun, Wei Liu, Weijia Shi, Huaxiu Yao, Linjun Zhang,
871 Andrew Y. Ng, James Zou, Sanmi Koyejo, Yejin Choi, Percy Liang, and Niklas Muennighoff. Uq:
872 Assessing language models on unsolved questions, 2025. URL <https://arxiv.org/abs/2508.17580>.
- 873 OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden
874 Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko,
875 Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally
876 Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich,
877 Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghor-
878 bani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao
879 Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi,
880 Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong
881 Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts,
882 Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David
883 Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong,
884 Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang,
885 Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred
886 von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace
887 Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin,
888 Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian
889 O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever,
890 Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng,
891 Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish,
892 Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan
893 Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singh, Karina Nguyen, Karl
894 Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin
895 Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus,
896 Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk,
897 Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko
898 Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz,
899 Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe,
900 Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang,
901 Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowd-
902 hury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg
903 Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias,
904 Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny
905 Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi
906 Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago
907 Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani
908 Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir
909 Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted
910 Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng,
911 Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie
912 Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou,
913 Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai,
914 Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. Openai o1 system card,
915 2024. URL <https://arxiv.org/abs/2412.16720>.
- 916 Yonatan Oren, Nicole Meister, Niladri S Chatterji, Faisal Ladhak, and Tatsunori Hashimoto. Proving
917 test set contamination in black-box language models. In *The Twelfth International Conference on*
918 *Learning Representations*, 2023.

- 920 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin
 921 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the
 922 finest text data at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.
- 923 Tomer Porian, Mitchell Wortsman, Jenia Jitsev, Ludwig Schmidt, and Yair Carmon. Resolving
 924 discrepancies in compute-optimal scaling of language models. *Advances in Neural Information
 925 Processing Systems*, 37:100535–100570, 2024.
- 926 Kun Qian, Shunji Wan, Claudia Tang, Youzhi Wang, Xuanming Zhang, Maximillian Chen, and Zhou
 927 Yu. Varbench: Robust language model benchmarking through dynamic variable perturbation, 2024.
 928 URL <https://arxiv.org/abs/2406.17681>.
- 929 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions
 930 for machine comprehension of text. In Jian Su, Kevin Duh, and Xavier Carreras (eds.), *Proceedings
 931 of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392,
 932 Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/
 933 D16-1264. URL <https://aclanthology.org/D16-1264/>.
- 934 Anka Reuel, Benjamin Bucknall, Stephen Casper, Timothy Fist, Lisa Soder, Onni Aarne, Lewis
 935 Hammond, Lujain Ibrahim, Alan Chan, Peter Wills, Markus Anderljung, Ben Garfinkel, Lennart
 936 Heim, Andrew Trask, Gabriel Mukobi, Rylan Schaeffer, Mauricio Baker, Sara Hooker, Irene
 937 Solaiman, Sasha Luccioni, Nitarshan Rajkumar, Nicolas Moës, Jeffrey Ladish, David Bau, Paul
 938 Bricman, Neel Guha, Jessica Newman, Yoshua Bengio, Tobin South, Alex Pentland, Sanmi Koyejo,
 939 Mykel Kochenderfer, and Robert Trager. Open problems in technical AI governance. *Transactions
 940 on Machine Learning Research*, 2025. ISSN 2835-8856. URL [https://openreview.net/
 941 forum?id=1n04qFMiS0](https://openreview.net/forum?id=1n04qFMiS0). Survey Certification.
- 942 Martin Riddell, Ansong Ni, and Arman Cohan. Quantifying contamination in evaluating code
 943 generation capabilities of language models. In Lun-Wei Ku, Andre Martins, and Vivek Sri Kumar
 944 (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics
 945 (Volume 1: Long Papers)*, pp. 14116–14137, Bangkok, Thailand, August 2024. Association for
 946 Computational Linguistics. doi: 10.18653/v1/2024.acl-long.761. URL <https://aclanthology.org/2024.acl-long.761/>.
- 948 Manley Roberts, Himanshu Thakur, Christine Herlihy, Colin White, and Samuel Dooley. To the
 949 cutoff... and beyond? a longitudinal perspective on LLM data contamination. In *The Twelfth
 950 International Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?id=m2NVG4Htxs](https://openreview.net/

 951 forum?id=m2NVG4Htxs).
- 952 Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. White-
 953 box vs black-box: Bayes optimal strategies for membership inference. In *International Conference
 954 on Machine Learning*, pp. 5558–5567. PMLR, 2019.
- 955 Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and
 956 Eneko Agirre. NLP evaluation in trouble: On the need to measure LLM data contamination for
 957 each benchmark. In *The 2023 Conference on Empirical Methods in Natural Language Processing*,
 958 2023. URL <https://openreview.net/forum?id=KivNpBsfAS>.
- 959 Oscar Sainz, Iker García-Ferrero, Alon Jacovi, Jon Ander Campos, Yanai Elazar, Eneko Agirre,
 960 Yoav Goldberg, Wei-Lin Chen, Jenny Chim, Leshem Choshen, Luca D’Amico-Wong, Melissa
 961 Dell, Run-Ze Fan, Shahriar Golchin, Yucheng Li, Pengfei Liu, Bhavish Pahwa, Ameya Prabhu,
 962 Suryansh Sharma, Emily Silcock, Kateryna Solonko, David Stap, Mihai Surdeanu, Yu-Min
 963 Tseng, Vishaal Udarrao, Zengzhi Wang, Ruijie Xu, and Jinglin Yang. Data contamination
 964 report from the 2024 CONDA shared task. In Oscar Sainz, Iker García Ferrero, Eneko Agirre,
 965 Jon Ander Campos, Alon Jacovi, Yanai Elazar, and Yoav Goldberg (eds.), *Proceedings of the
 966 1st Workshop on Data Contamination (CONDA)*, pp. 41–56, Bangkok, Thailand, August 2024.
 967 Association for Computational Linguistics. doi: 10.18653/v1/2024.conda-1.4. URL <https://aclanthology.org/2024.conda-1.4/>.
- 969 Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes.
 970 ML-leaks: Model and data independent membership inference attacks and defenses on machine
 971 learning models. *arXiv preprint arXiv:1806.01246*, 2018.

- 972 Rylan Schaeffer. Pretraining on the test set is all you need, 2023. URL <https://arxiv.org/abs/2309.08632>.
- 974 Rylan Schaeffer, Zachary Robertson, Akhilan Boopathy, Mikail Khona, Kateryna Pistunova, Ja-
975 son William Rocks, Ila R Fiete, Andrey Gromov, and Sanmi Koyejo. Double descent demystified:
976 Identifying, interpreting & ablating the sources of a deep learning puzzle. In *The Third Blogpost
977 Track at ICLR 2024*.
- 978 Rylan Schaeffer, Joshua Kazdan, and Yegor Denisov-Blanch. Min-p, max exaggeration: A critical
979 analysis of min-p sampling in language models, 2025a. URL <https://arxiv.org/abs/2506.13681>.
- 981 Rylan Schaeffer, Hailey Schoelkopf, Brando Miranda, Gabriel Mukobi, Varun Madan, Adam Ibrahim,
982 Herbie Bradley, Stella Biderman, and Sanmi Koyejo. Why has predicting downstream capabilities
983 of frontier AI models with scale remained elusive? In *Forty-second International Conference on
984 Machine Learning*, 2025b. URL <https://openreview.net/forum?id=I1NtlLvJa1>.
- 985 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen,
986 and Luke Zettlemoyer. Detecting pretraining data from large language models. In *The Twelfth
987 International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=zWqr3MQuNs>.
- 989 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks
990 against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18.
991 IEEE, 2017.
- 992 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and
993 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.
994 In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard (eds.),
995 *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp.
996 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.
997 URL <https://aclanthology.org/D13-1170/>.
- 998 Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization
999 without overfitting: Analyzing the training dynamics of large language models. *Advances in Neural
1000 Information Processing Systems*, 35:38274–38290, 2022.
- 1001 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
1002 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-
1003 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
1004 Wenjin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
1005 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
1006 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
1007 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
1008 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
1009 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
1010 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
1011 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
1012 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
1013 2023. URL <https://arxiv.org/abs/2307.09288>.
- 1014 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
1015 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing
1016 systems*, 30, 2017.
- 1017 Xinyi Wang, Antonis Antoniades, Yanai Elazar, Alfonso Amayuelas, Alon Albalak, Kexun Zhang,
1018 and William Yang Wang. Generalization v.s. memorization: Tracing language models’ capabilities
1019 back to pretraining data. In *The Thirteenth International Conference on Learning Representations*,
1020 2025. URL <https://openreview.net/forum?id=IQxBDLmVpt>.
- 1021 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
1022 Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International
1023 Conference on Learning Representations*, 2022.

- 1024 Chunqiu Steven Xia, Yinlin Deng, and Lingming Zhang. Top leaderboard ranking= top coding profi-
1025 ciency, always? evoeval: Evolving coding benchmarks via llm. *arXiv preprint arXiv:2403.19114*,
1026 2024.
- 1027 Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih
1028 Ghazi, and Ravi Kumar. On memorization of large language models in logical reasoning, 2025.
1029 URL <https://arxiv.org/abs/2410.23123>.
- 1030 Cheng Xu, Shuhao Guan, Derek Greene, and M-Tahar Kechadi. Benchmark data contamination of
1031 large language models: A survey, 2024a. URL <https://arxiv.org/abs/2406.04244>.
- 1032 Fengli Xu, Qianyue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong
1033 Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, Chenyang Shao, Yuwei Yan, Qinglong Yang,
1034 Yiwen Song, Sijian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao, and Yong Li. Towards large
1035 reasoning models: A survey of reinforced reasoning with large language models, 2025. URL
1036 <https://arxiv.org/abs/2501.09686>.
- 1037 Ruijie Xu, Zengzhi Wang, Run-Ze Fan, and Pengfei Liu. Benchmarking benchmark leakage in large
1038 language models, 2024b. URL <https://arxiv.org/abs/2404.18824>.
- 1039 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
1040 Gao, Chengan Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
1041 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
1042 Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
1043 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
1044 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
1045 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
1046 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
1047 Qiu. Qwen3 technical report, 2025a. URL <https://arxiv.org/abs/2505.09388>.
- 1048 Shuo Yang, Wei-Lin Chiang, Lianmin Zheng, Joseph E Gonzalez, and Ion Stoica. Rethinking
1049 benchmark and contamination for language models with rephrased samples. *arXiv preprint*
1050 *arXiv:2311.04850*, 2023.
- 1051 Zhen Yang, Hongyi Lin, Yifan He, Jie Xu, Zeyu Sun, Shuo Liu, Pengpeng Wang, Zhongxing Yu, and
1052 Qingyuan Liang. Rethinking the effects of data contamination in code intelligence, 2025b. URL
1053 <https://arxiv.org/abs/2506.02791>.
- 1054 Feng Yao, Yufan Zhuang, Zihao Sun, Sunan Xu, Animesh Kumar, and Jingbo Shang. Data contami-
1055 nation can cross language barriers. In *Proceedings of the 2024 Conference on Empirical Methods*
1056 *in Natural Language Processing*, pp. 17864–17875, 2024.
- 1057 Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning:
1058 Analyzing the connection to overfitting. In *2018 IEEE 31st computer security foundations*
1059 *symposium (CSF)*, pp. 268–282. IEEE, 2018.
- 1060 Sajjad Zarifzadeh, Philippe Liu, and Reza Shokri. Low-cost high-power membership inference
1061 attacks, 2023.
- 1062 Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav
1063 Raja, Charlotte Zhuang, Dylan Slack, Qin Lyu, Sean Hendryx, Russell Kaplan, Michele Lunati,
1064 and Summer Yue. A careful examination of large language model performance on grade school
1065 arithmetic, 2024a. URL <https://arxiv.org/abs/2405.00332>.
- 1066 Jie Zhang, Debeshee Das, Gautam Kamath, and Florian Tramèr. Membership inference attacks
1067 cannot prove that a model was trained on your data. *arXiv preprint arXiv:2409.19798*, 2024b.
- 1068 Zhe Zhang, Runlin Liu, Aishan Liu, Xingyu Liu, Xiang Gao, and Hailong Sun. Dynamic benchmark
1069 construction for evaluating large language models on real-world codes, 2025. URL <https://arxiv.org/abs/2508.07180>.
- 1071 Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin,
1072 Ji-Rong Wen, and Jiawei Han. Don't make your llm an evaluation benchmark cheater, 2023. URL
1073 <https://arxiv.org/abs/2311.01964>.

1074 **A Related Work**

1075 **Data Contamination and its Consequences** Test set contamination, where benchmark data is
1076 included in pretraining corpora, is widely recognized as a threat to valid model evaluation, as it can
1077 lead to inflated performance metrics. Numerous survey and position papers have documented the
1078 various ways contamination can occur and have called for routine audits and transparent reporting for
1079 all benchmarks (Sainz et al., 2023, 2024; Deng et al., 2024a; Xu et al., 2024a; Reuel et al., 2025).
1080 Empirical studies of large web-scale datasets have confirmed significant overlap and duplication
1081 between training and test sets (Dodge et al., 2021). Research focused on ensuring benchmark integrity
1082 has identified multiple ways that language models might "cheat" on evaluations if contamination
1083 is not properly managed (Zhou et al., 2023; Dong et al., 2024). For instance, analyses of popular
1084 mathematics benchmarks have revealed signals of data leakage and potential overfitting (Zhang
1085 et al., 2024a). Ongoing community efforts and open-source audits continue to measure the extent
1086 of contamination across different models and datasets (Li et al., 2024). The risks extend beyond
1087 evaluation integrity; scaling studies indicate that poisoning risks increase with model size, as larger
1088 models learn harmful behaviors from minuscule amounts of poisoned data far more rapidly than
1089 smaller models, underscoring the necessity of robust data curation (Bowen et al., 2025). As a
1090 cautionary illustration, Schaeffer (2023) demonstrated that pretraining on the test set is a trivial
1091 path to strong benchmark performance, reinforcing the importance of rigorous decontamination and
1092 auditing.

1093 **Controlled Contamination During Pretraining** A line of research directly investigates the causal
1094 effects of contamination by intentionally adding benchmark data to pretraining corpora and observing
1095 the results. Magar & Schwartz (2022) interleaved task-specific datasets into a general text corpus
1096 during pretraining, varying the duplication rate of the leaked examples. They differentiated between
1097 "memorization" (storing examples) and "exploitation" (using stored examples to boost test scores),
1098 finding that both model size and the number of repetitions increased exploitation. Jiang et al. (2024)
1099 pretrained models from scratch on corpora containing either only the inputs ("text-only") or the full
1100 input-output pairs ("ground-truth") of benchmark examples, sweeping the contamination frequency.
1101 They observed significant performance gains when ground-truth pairs were used and showed that
1102 simple n-gram-based detection methods could be bypassed by paraphrasing or partial data leaks.
1103 The problem also transcends language barriers; Yao et al. (2024) demonstrated a cross-lingual
1104 contamination channel where continuing to pretrain a model on non-English translations of English
1105 benchmarks led to material improvements on the original English tests, a form of contamination that
1106 string-matching would not detect. At a larger scale, Bordt et al. (2025) varied the repetition count of
1107 leaked examples, model size (up to 1.6B parameters), and the total training token budget, finding
1108 that performance scales predictably with size and repetition. They also showed that sufficiently long
1109 training on abundant unique data could mitigate or even reverse the effects of earlier contamination.
1110 In the context of machine translation, Kocyigit et al. (2025) injected source-target pairs into the
1111 pretraining data of 1B and 8B parameter models, quantifying significant overestimation in BLEU
1112 scores, with larger models and low-resource languages showing more pronounced effects. Together,
1113 these causal intervention studies provide clear evidence that language models memorize and leverage
1114 benchmark data when it is present during pretraining.

1115 **Repeated Data and Memorization Dynamics** Closely related is the study of memorization
1116 dynamics, particularly how repeated data affects model behavior. Hernandez et al. (2022) trained
1117 models where a small portion of the data was repeated many times, observing strong double descent
1118 phenomena (Advani et al., 2020; Belkin et al., 2019; Adlam & Pennington, 2020; Bordelon et al.,
1119 2020; Schaeffer et al.) and showing that repeating just 0.1% of tokens 100 times could significantly
1120 degrade generalization. Studies tracking exact-sequence memorization have shown that larger
1121 models not only memorize more content and at a faster rate but also forget less over the course of
1122 training (Tirumala et al., 2022). Carlini et al. (2023) quantified log-linear relationships between
1123 verbatim generation and model size, data duplication count, and prompt length. Other work has
1124 explored the feasibility of *forecasting* whether a model will memorize a specific string, finding that
1125 accurate prediction is possible but may require a substantial portion of the target model's pretraining
1126 compute (Biderman et al., 2023). Beyond explicit repetition, Duan et al. (2025) discovered *latent
1127 memorization*, where memorized sequences that are not obvious at a final checkpoint can persist
1128 and be revealed later, posing privacy risks. Finally, memorization appears to be task-dependent:
1129 Wang et al. (2025) observed stronger memorization for knowledge-intensive QA, whereas machine

1130 translation and mathematical reasoning demonstrated greater novelty. Memorization also interacts
1131 with logical reasoning; using dynamically generated puzzles, Xie et al. (2025) showed that models
1132 could be fine-tuned to perfectly memorize training examples yet failed on slight variations, even as
1133 their genuine reasoning abilities also improved, revealing a complex balance between the two.

1134 **Detecting and Proving Contamination** Another significant area of research focuses on detecting
1135 or proving test set contamination in existing models. Oren et al. (2023) and Ni et al. (2025) proposed
1136 statistical tests with provable control over false positives by testing if a benchmark's canonical
1137 ordering is statistically privileged over random shuffles. Shi et al. (2024) introduced Min- $k\%$ -Prob
1138 to determine if a sequence likely appeared in pretraining using only black-box probabilities. Two
1139 related works from Golchin & Surdeanu (2023, 2024) frame detection as a multiple-choice "quiz"
1140 and use temporal information about model training windows versus benchmark release dates, a
1141 strategy also used by Roberts et al. (2024). Broader audits have aimed to quantify leakage and
1142 decontamination across a wide range of tasks and models (Xu et al., 2024b; Deng et al., 2024b; Li
1143 et al., 2024), while Yang et al. (2023) showed that rephrasing benchmark questions can often bypass
1144 n-gram filters. In the domain of code generation, Riddell et al. (2024) quantified contamination
1145 in popular coding benchmarks and connected the degree of overlap to performance differences.
1146 Matton et al. (2024) cataloged various channels for leakage and released a dataset (LBPP) to help
1147 mitigate these issues. Complementing these audits, Yang et al. (2025b) systematically tested fine-
1148 grained contamination scenarios in code intelligence across different model types, finding that paired
1149 contamination substantially affects LLMs under a pretraining-plus-inference paradigm but has limited
1150 effect under a pretrain–finetune–inference pipeline. Other work has also provided instruments for
1151 detecting the origins of chain-of-thought sequences (Li et al., 2025).

1152 **Preventing Test Set Contamination** The growing concern over contamination has spurred the de-
1153 velopment of new methods for creating benchmarks. These include dynamically updated benchmarks
1154 (Jain et al., 2025; Xia et al., 2024; Zhang et al., 2025; Qian et al., 2024) and private or restricted-
1155 access benchmarks (Zhang et al., 2024a; Glazer et al., 2025). Recently, Nie et al. (2025) released a
1156 benchmark consisting of unsolved scientific questions, which, by its nature, prevents models from
1157 being trained on the correct solutions.

1158 **Retrieval- and Agent-Time Contamination** As model evaluation evolves from static prompting
1159 to using tool-augmented agents, the risk of contamination expands. Han et al. (2025) introduced
1160 search-time contamination, where an agent retrieves benchmark questions and answers from the web
1161 during its evaluation process, which can artificially inflate its performance.

1162 **Membership Inference Attacks** The field of Membership Inference Attacks (MIA) aims to
1163 determine if a specific data point was used to train a model, given only access to the model itself
1164 (Shokri et al., 2017). This is highly relevant to contamination, as detection can be viewed as an
1165 MIA problem. While the MIA literature is extensive in computer vision (Yeom et al., 2018; Salem
1166 et al., 2018; Sablayrolles et al., 2019; Jagielski et al., 2024), it has more recently been applied to
1167 language models (Carlini et al., 2021; Zarifzadeh et al., 2023; Shi et al., 2024; Mattern et al., 2023;
1168 Li et al., 2023). However, progress in sequence-level MIA for language models has been complicated
1169 by issues such as flawed evaluations (Meeus et al., 2024; Zhang et al., 2024b; Jiang et al., 2025).
1170 Duan et al. (2024) argue that membership can be inherently "blurry" for natural language. Das et al.
1171 (2024) and Meeus et al. (2024) report that existing MIA testbeds suffer from distribution shifts. Kong
1172 et al. (2023) refute MIAs with a theoretical attack, and Liu et al. (2025) and Mangaokar et al. (2025)
1173 demonstrate fundamental limitations and exploits of n-gram based methods. Due to these challenges,
1174 recent work explores strengthening the membership signal by using multiple correlated sequences as
1175 input (Maini et al., 2021; Kandpal et al., 2023; Maini et al., 2024), which aligns more closely with
1176 detecting contamination of an entire test set rather than a single example (Golchin & Surdeanu, 2023;
1177 Oren et al., 2023).