

GHAJAR EXHIBIT 50

1 COOLEY LLP
BOBBY GHAJAR (198719)
2 (bghajar@cooley.com)
COLETTE GHAZARIAN (322235)
3 (cghazarian@cooley.com)
1333 2nd Street, Suite 400
4 Santa Monica, California 90401
Telephone: (310) 883-6400

5 COOLEY LLP
6 MARK WEINSTEIN (193043)
(mweinstein@cooley.com)
7 KATHLEEN HARTNETT (314267)
(khartnett@cooley.com)
8 JUDD LAUTER (290945)
(jlauter@cooley.com)
9 ELIZABETH L. STAMESHKIN (260865)
(lstameshkin@cooley.com)
10 3175 Hanover Street
Palo Alto, CA 94304-1130
11 Telephone: (650) 843-5000

12 CLEARY GOTTlieb STEEN & HAMILTON LLP
ANGELA L. DUNNING (212047)
13 (adunning@cgsh.com)
1841 Page Mill Road, Suite 250
14 Palo Alto, CA 94304
Telephone: (650) 815-4131

15 *[Full Listing on Signature Page]*

16 *Counsel for Defendant Meta Platforms, Inc.*

17
18 **UNITED STATES DISTRICT COURT**
19 **NORTHERN DISTRICT OF CALIFORNIA**

20 RICHARD KADREY, et al.,
21 Individual and Representative Plaintiffs,
22 v.
23 META PLATFORMS, INC., a Delaware
24 corporation;
25 Defendant.

Case No. 3:23-cv-03417-VC

**DEFENDANT META PLATFORMS, INC.'S
FURTHER SUPPLEMENTAL AND AMENDED
RESPONSES AND OBJECTIONS TO
PLAINTIFFS' FIRST SET OF
INTERROGATORIES**

Trial Date: None
Date Action Filed: July 7, 2023

Interrogatory, Meta will conduct a reasonable, proportionate search for non-privileged, relevant, responsive information within its possession, custody, or control.

17. In responding to all Interrogatories, Meta will comply with the requirements of the Federal Rules of Evidence and Federal Rule of Civil Procedure 26.

III. OBJECTIONS AND RESPONSES TO INDIVIDUAL INTERROGATORIES

INTERROGATORY NO. 1:

Describe in detail the data You have used to train or otherwise develop the Meta Language Models, Including, for each:

- a. How You obtained the data, e.g., by scraping the data, purchasing it from third parties, or by other means;
- b. All sources of Data, including any third parties that provided data sets;
- c. To the extent the data was derived from publicly available websites, a list of all such websites and, for each, the percentage of the data corpus that is derived from that website;
- d. The categories of content included in the data and the extent to which each category is represented in the data corpus (i.e., as a percentage of data used to train the model);
- e. All policies and procedures Related to identifying, assessing, vetting and selecting sources of data for the model.

RESPONSE TO INTERROGATORY NO. 1:

Meta incorporates by reference its objections and definitions above, including to the terms “You” and “Meta Language Models.” Meta further notes that the capitalized term “Related” is not defined; Meta construes that term coextensively with “concerning.”

As an initial matter, Meta objects to this Interrogatory because it consists of multiple, separate Interrogatories, each which count toward Plaintiffs’ limit under the Federal Rules. For example, the question about what data used to train a model is separate from how it was obtained, and further, subparts (d) and (e) are not subsumed within and necessarily related to the primary question, and purport to require a calculation of percentages of data, and separate identification of “policies” and “procedures” for (1) identifying, (2) assessing, (3) vetting, and (4) selecting data. This Interrogatory consists of *at least* three Interrogatories, and depending on how it is interpreted,

1 many more. In answering the Interrogatory, Meta does not waive this objection.

2 Meta objects to this Interrogatory because, on its face, it does not exclude legal advice or
3 opinions, which are subject to attorney-client privilege and/or attorney work product doctrine, in
4 particular as to subpart (e). Meta will not produce privileged materials or attorney work product.

5 Meta objects to this Interrogatory as vague and ambiguous as to the term “data,” which is
6 undefined. Meta will construe “data” to mean Training Data (as construed above).

7 Meta objects to this Interrogatory as vague, ambiguous, and unintelligible as to “percentage
8 of that data corpus that is derived from that website” because “data corpus” is undefined, and Meta
9 is accordingly unable to interpret and respond to subpart (c). Even if “data corpus” were defined,
10 the subject matter of subpart (c) would be overbroad, unduly burdensome, and disproportionate to
11 the needs of the case and seeks information that is not relevant to the parties’ claims and defenses.
12 Meta will not respond to subpart (c).

13 Meta objects to the undefined phrase “categories of content, which is vague, ambiguous,
14 and unintelligible.

15 Meta objects to this Interrogatory to the extent that it seeks information that is not within
16 Meta’s possession, custody, or control.

17 Subject to and without waiving the foregoing objections, and pursuant to the terms of the
18 Protective Order and the ESI Order, Meta responds as follows: Meta incorporates by reference the
19 identification of datasets used to train Llama 1 that is included in the publicly available paper
20 “LLaMA: Open and Efficient Foundation Language Models.” Such datasets were used to train
21 Llama 1. Meta will produce a copy of that paper in its forthcoming production pursuant to Rule
22 33(d).

23 Meta will conduct a reasonable search for additional non-privileged information or, in
24 accordance with Rule 33(d), documents in Meta’s possession, custody, or control, sufficient to
25 show any other datasets used to train the Meta Language Models (as construed above), as well as
26 policies and procedures for identifying, assessing, vetting, and selecting sources of data for those
27 models.

28 Discovery is ongoing and Meta will also supplement its response to this Interrogatory to

identify the sources of such datasets and general categories of data within them, to the extent that such information is within Meta's possession, custody, or control.

Discovery is continuing and Meta reserves the right to supplement or amend its response at a later time.

Meta's First Supplemental and Amended Response to Interrogatory No. 1:

Subject to and without waiving the foregoing objections, and pursuant to the terms of the Protective Order, Meta responds as follows.

This response is designated as Highly Confidential – Attorney's Eyes Only under the Protective Order.

Based on its reasonable investigation, Meta identifies the following datasets as containing material used to train the Llama Models (as construed above), including pretraining and/or finetuning, as well the locations from which Meta believes they were obtained:

<u>Dataset name</u>	<u>Llama 1</u>	<u>Llama 2</u>	<u>Llama 3</u>	<u>URL or Other Location</u>
Stack Exchange	Yes	Yes	Yes	https://archive.org/details/stackexchange
books3	Yes	Yes	Yes	https://the-eye.eu/public/AI/pile_preliminary_components/books3.tar.gz
Project Gutenberg	Yes	Yes	Yes	https://www.gutenberg.org
Arxiv	Yes	Yes	Yes	https://www.arxiv.org
Github	Yes	Yes	Yes	https://www.github.com
C4	Yes	Yes	Yes	https://www.tensorflow.org/datasets/catalog/c4
CCNet	Yes	Yes	No	https://github.com/speedinghzl/CCNet/blob/master/LICENSE
CC-stories	Yes	Yes	No	https://github.com/tensorflow/models/tree/archive/research/lm_commonsense#1-download-data-files
The Stack	Yes	Yes	Yes	https://huggingface.co/datasets/bigcode/the-stack
Wikipedia	Yes	Yes	Yes	https://en.wikipedia.org/wiki/Wikipedia:Database_download
	No	Yes	Yes	
	No	Yes	No	
	No	Yes	Yes	

[illegible]

1		No	No	Yes	
2		No	No	Yes	
3		No	No	Yes	
4					
5					
6					
7					
8	Libgen	No	No	Yes	https://libgen.is
9		No	No	Yes	
10		No	No	Yes	
11		No	No	Yes	
12		No	No	Yes	
13		No	No	Yes	
14		No	No	Yes	
15		No	No	Yes	
16		No	No	Yes	
17		No	No	Yes	
18		No	No	Yes	
19		No	No	Yes	
20		No	No	Yes	
21		No	No	Yes	
22		No	No	Yes	
23					
24		No	No	Yes	
25		No	No	Yes	
26		No	No	Yes	
27		No	No	Yes	
28					

1		No	No	Yes	
2		No	No	Yes	
3					
4		No	No	Yes	
5					
6		No	No	Yes	
7		No	No	Yes	
8					
9					
10		No	No	Yes	
11		No	No	Yes	
12					
13					
14					
15					
16					
17					
18		No	No	Yes	
19		No	No	Yes	
20		No	No	Yes	
21		No	No	Yes	
22		No	No	Yes	
23		No	No	Yes	
24		No	No	Yes	
25		No	No	Yes	
26		No	No	Yes	
27		No	No	Yes	
28		No	No	Yes	

1 The annotations data used for finetuning Llama 2 are identified in Table 6 of the paper titled
2 “Llama 2: Open Foundation and Fine-Tuned Chat Models.” Except for the “Meta (Safety &
3 Helpfulness)” data, which was obtained from Meta’s vendors, namely, [REDACTED] these annotations
4 datasets were sourced from publicly available sources, such as Github and Hugging Face. Meta has
5 also entered agreements with [REDACTED] to provide annotations.

6 In addition, for Llama 3.1, Meta used publicly available data sourced from [REDACTED]
7 [REDACTED] to train the model, as well as a variety of synthetic data.

8 The process for selecting datasets for use in pre-training of the Meta Language Models (as
9 construed above) was informed by what data was available, whether the development team believed
10 that the data would help the model achieve optimal results against industry benchmarks, and PXFN
11 review, *i.e.*, cross-functional review by legal, privacy, and/or policy personnel. Each of the above
12 external datasets was required to undergo PXFN review prior to training of the Meta Language
13 Models (as construed above). Any issues related to intellectual property are regarded as legal in
14 nature. Review and consideration of those issues is therefore the responsibility of Meta’s legal
15 team, rather than Meta privacy or policy personnel, and is subject to attorney-client privilege and/or
16 work product doctrine.

17 From the development team’s perspective, decisions around which datasets to use for Llama
18 1 were influenced by the development of other large language models, in particular DeepMind’s
19 Chinchilla and the corresponding paper “Training Compute-Optimal Large Language Models.” At
20 the time, researchers regarded DeepMind’s Chinchilla as state of the art, and the team developing
21 the first version of Llama was motivated to reproduce Chinchilla’s results on industry benchmarks
22 (e.g., MMLU, BoolQ, PIQA, etc.) using their own model architecture. Using the same or similar
23 dataset diversity allowed the team to better compare the effectiveness of the respective models.
24 Llama 2 was largely trained on the same datasets as Llama 1.

25 With respect to Llama 3, whether a particular dataset was used in the training of the model
26 was driven by a number of considerations, including:

- 27 • Dataset size – It is understood that LLMs, such as the Meta Language Models (as construed
28 above), require large volumes of text data in order to achieve high performance across

industry benchmarks. In general, the more data the models train on, the better the performance of the model. That is, there is a rough correlation between the number of unique token strings within a dataset and the performance of the models on downstream tasks, such as the ability to answer questions. Accordingly, larger datasets are preferred to smaller datasets.

- Dataset diversity – Datasets with greater diversity of subject matter, a variety of lengths and human/computer languages, and different styles of writing or conversation help enable the models to be more flexible and adaptable to different contexts.
- Dataset quality – Related to the diversity of the dataset is the extent to which undesirable data (such as repetitive data, factually incorrect data, or harmful or toxic data) can be filtered from the dataset without degrading dataset usefulness.

The data mix that will achieve the best results against benchmarks (e.g., MMLU, GSM8K, BoolQ, PIQA, CommonsenseQA, etc.) is difficult to determine in advance. Accordingly, the Meta Language Model (as construed above) development teams performed small scale experiments prior to full scale pre-training to evaluate optimal data mix proportions. Pursuant to Rule 33(d), Meta also refers Plaintiffs to the paper titled “The Llama 3 Herd of Models,” published by Meta on July 23, 2024, for further information.

INTERROGATORY NO. 3:

Describe in detail the RLHF process for each Meta Language Model. Include in Your response:

- a. Examples of types of experts who write questions and answers for use in RLHF;
- b. Examples of questions and answers;
- c. An explanation of the rating system or method of evaluation for the Meta Language Model’s responses;
- d. A description of the RLHF You actually undertook in order to correct or remediate any Meta Language Model’s propensity to emit protected expression from its Training Data.

RESPONSE TO INTERROGATORY NO. 3:

Meta incorporates by reference its objections and definitions above, including to the terms

1 Dated: December 13, 2024

COOLEY LLP

2 By: /s/ Judd Lauter

3 Bobby Ghajar
4 Mark Weinstein
5 Kathleen Hartnett
6 Phillip Morton
7 Judd Lauter
8 Elizabeth L. Stameshkin
9 Matthew Brigham
10 Colette Ghazarian
11 Juan Pablo Gonzalez
12 Cole A. Poppell

LEX LUMINA PLLC
Mark A. Lemley

CLEARY GOTTlieb STEEN &
HAMILTON LLP
Angela L. Dunning

Attorneys for Defendant
META PLATFORMS, INC.

13 *Full Counsel List*

14 COOLEY LLP
15 PHILLIP MORTON (*pro hac vice*)
16 (pmorton@cooley.com)
17 COLE A. POPPELL (*pro hac vice*)
18 (cpoppell@cooley.com)
1299 Pennsylvania Avenue, NW, Suite 700
Washington, DC 20004-2400
Telephone: (202) 842-7800

19 COOLEY LLP
20 MATTHEW BRIGHAM (191428)
21 (mbrigham@cooley.com)
22 JUAN PABLO GONZALEZ (334470)
23 (jgonzalez@cooley.com)
3175 Hanover Street
Palo Alto, CA 94304-1130
Telephone: (650) 843-5000

24 LEX LUMINA PLLC
25 MARK A. LEMLEY (155830)
26 (mlemley@lex-lumina.com)
745 Fifth Avenue, Suite 500
New York, NY 10151
Telephone: (646) 898-2055

VERIFICATION

I, Michael Clark, declare:

I am an employee of Meta Platforms, Inc. ("Meta"), a corporation organized and existing under the laws of Delaware, which is the Defendant in the above-entitled action, and I have been authorized to make this verification on its behalf.

I have read the following documents:

- Meta's Further Supplemental and Amended Responses and Objections to Plaintiffs' First Set of Interrogatories.
- Meta's Further Supplemental and Amended Responses and Objections to Plaintiffs' Second Set of Interrogatories.
- Meta's First Supplemental Responses and Objections to Plaintiffs' Third Set of Interrogatories

I believe, based on personal knowledge or upon information and belief, that those responses are true and correct.

I declare under penalty of perjury under the laws of the United States that the foregoing is true and correct.

Executed at Denver, Colorado on December 13, 2024.


Michael Clark