

1 QUINN EMANUEL URQUHART &
SULLIVAN, LLP
2 Sean S. Pak (SBN 219032)
seanpak@quinnemanuel.com
3 50 California Street, 22nd Floor
San Francisco, CA 94111
4 Telephone: (415) 875-6600
Facsimile: (415) 875-6700
5

Andrew H. Schapiro (admitted *pro hac vice*)
6 andrewschapiro@quinnemanuel.com
191 N. Wacker Drive, Suite 2700
7 Chicago, Illinois 60606
Telephone: (312) 705-7400
8 Facsimile: (312) 705-4001

9 Alex Spiro (admitted *pro hac vice*)
alexspiro@quinnemanuel.com
10 295 Fifth Avenue
New York, NY 10016
11 Telephone: (212) 849-7000
Facsimile: (212) 849-7100
12

Rachael L. McCracken (SBN 252660)
13 rachaelmccracken@quinnemanuel.com
865 South Figueroa Street, 10th Floor
14 Los Angeles, CA 90017
Telephone: (213) 443-3000
15 Facsimile: (213) 443-3100

16 *Attorneys for Defendant,*
NVIDIA Corporation

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

20 ABDI NAZEMIAN, BRIAN KEENE,
21 STEWART O’NAN, ANDRE DUBUS III and
SUSAN ORLEAN, individually and on behalf
22 of all others similarly situated,

23 Plaintiffs,

24 v.

25 NVIDIA CORPORATION,

26 Defendant.

Master File Case No. 4:24-cv-01454-JST (SK)

**DEFENDANT’S RESPONSE TO MOTION
FOR RELIEF FROM NONDISPOSITIVE
PRETRIAL ORDER OF MAGISTRATE
JUDGE**

27
28

1 **PRELIMINARY STATEMENT**

2 Plaintiffs ask this Court to overturn the Magistrate Judge’s decision which correctly limited
3 discovery to the dataset on which Plaintiffs based their infringement claim—consistent with the
4 decisions of every court to examine the issue. This case is about whether NVIDIA’s use of the
5 Books3 dataset to train its “NeMo Megatron” family of large language models (“LLMs”) is
6 copyright infringement or fair use. Plaintiffs focused their case on a single dataset (Books3)—which
7 is a subset of a larger dataset called “The Pile.” ECF 1 (“Compl.”) ¶¶ 25-26, 29, 31 (identifying only
8 Books3 as containing Plaintiffs’ works). Plaintiffs did *not* allege that their works are found in any
9 dataset other than Books3. Similarly, Plaintiffs alleged only that NVIDIA infringed their copyrights
10 by using Books3 to train a single model family, the NeMo Megatron family of LLMs. *Id.* ¶¶ 23, 31,
11 34 (“To train the NeMo Megatron language models, NVIDIA copied The Pile dataset.”).

12 Despite these narrow allegations, Plaintiffs have sought sweeping discovery into LLMs and
13 datasets (which Plaintiffs call “shadow libraries”) not identified in the Complaint and that Plaintiffs
14 did not allege were trained on or contain copies of their works. ECF 162. Plaintiffs argued to the
15 Magistrate Judge that discovery into models and datasets “beyond those named in the Complaint”
16 might allow them to make other claims than what they pled in the Complaint. *Id.* at 1-4. The
17 Magistrate Judge considered and correctly rejected these arguments, properly limiting discovery to
18 (1) “the dataset that Plaintiffs know that Defendant used and that contains Plaintiffs’ copyrighted
19 books”—the Pile—and (2) LLMs “in the Nemo Megatron family that exist and that trained on the
20 Pile dataset.” ECF 168 (the “Order”) at 2-3. Plaintiffs challenge only the decision on datasets.

21 Plaintiffs seek to overturn the Magistrate Judge’s decision and compel NVIDIA to run
22 Search Term 1 to fish for information untethered to the allegations in their Complaint. They
23 acknowledge as much in a footnote. *See* Mot. at 1 n.2 (“Should discovery show Nvidia used other
24 shadow libraries to train its LLMs, Plaintiffs may seek to expand the scope of models through a
25 discovery motion or leave to amend.”). Every court to examine requests like Plaintiffs’ for discovery
26 into datasets untethered to the operative complaint or the plaintiffs’ works has rightly denied those
27 requests, as the Order did. Because the Magistrate Judge committed no error, let alone clear error,
28 the Court should deny the Motion.

ARGUMENT

1
2 A magistrate judge’s non-dispositive order may only be overturned if there is “a definite and
3 firm conviction that a mistake has been committed.” *Perry v. Schwarzenegger*, 268 F.R.D. 344, 348
4 (N.D. Cal. 2010). “The court should not disturb the magistrate’s relevance determination except
5 where it is based on an erroneous conclusion of law or where the record contains no evidence on
6 which the magistrate rationally could have based that decision.” *Id.* (citation omitted). Plaintiffs here
7 fail to show the Order is “clearly erroneous or contrary to law.” Fed. R. Civ. P. 72(a).

8 ***Plaintiffs’ Direct Infringement Claim Is Limited To The Pile.*** The Complaint alleges
9 “NVIDIA copied The Pile dataset,” which “includes the Books3 dataset, which includes the
10 Infringed Works,” and that NVIDIA “made multiple copies of the Books3 dataset while training the
11 NeMo Megatron models.” Compl. ¶ 34. Nowhere do Plaintiffs allege their asserted works can be
12 found in any dataset other than The Pile or that NVIDIA used any dataset other than The Pile to
13 train the NeMo Megatron LLMs.¹ The Complaint’s sole generic reference to “shadow libraries”
14 does not allege NVIDIA used them or that they contain Plaintiffs’ works and does not purport to
15 base any cause of action on them. *Id.* ¶ 27 (noting “shadow libraries” “have long been of interest to
16 the AI-training community because they host and distribute vast quantities of unlicensed
17 copyrighted material”). This passing reference, without more, cannot make these datasets relevant
18 to Plaintiffs’ infringement claim. Further, Plaintiffs’ suggestion that NVIDIA did not deny having
19 “knowledge (or possession)” of these “shadow libraries” in its answer is misleading, at best. Mot.
20 at 3. The Complaint did not allege any such knowledge or possession, so it was not addressed in the
21 answer. Compl. ¶ 27.

22 ***There is no precedent for the untethered discovery Plaintiffs seek into unpled datasets.***
23 Like the Magistrate Judge’s Order, every case that has considered discovery requests for datasets
24 not identified in the complaints’ infringement allegations has rejected them. In *Anderson v. Stability*

25 _____
26 ¹ Shortly before the hearing before the Magistrate Judge, Plaintiffs asserted for the first time that
27 their works can be found in LibGen, Anna’s Archive, and Z-Library. Plaintiffs raise the same
28 argument in a single sentence and footnote. Mot. at 2 n.3 (citing a URL to a third-party website).
Plaintiffs do not dispute that the Complaint lacks allegations regarding these datasets, so their new
assertion is irrelevant and should be disregarded. In any event, Plaintiffs’ lone citation to a third-
party website, without any allegations directed at **NVIDIA**, cannot justify this new, unpled theory.

1 *AI Ltd.*, the plaintiffs’ complaint asserted infringement of works “contained in the LAION datasets,”
2 but plaintiffs nonetheless sought discovery into all datasets used to train the accused AI models. No.
3 23-cv-00201-WHO, ECF 307 at 3-5, ECF 314 at 2. The court denied the plaintiffs’ broad request,
4 finding that (1) their “expansive view of their case [wa]s at odds with the scope of the operative
5 complaint,” and (2) discovery into non-LAION datasets—which were not alleged to include
6 Plaintiffs’ copyrighted works—was a “fishing expedition” “disproportionate ‘to the needs of the
7 case’ and contrary to the limits imposed by Rule 26(b)(1).” *Id.* at 3-5 (citation omitted).

8 Similarly, in *In re Mosaic LLM*, the court granted a motion to dismiss and denied discovery
9 on similar facts, even where Plaintiffs later attempted to add allegations directed to additional
10 datasets, finding “Plaintiffs do not allege facts that could establish that the [AI] models are actually
11 trained on any shadow library websites, let alone those that contain Plaintiffs’ works.” 2025 WL
12 2402677, at *2-3 (N.D. Cal. Aug. 19, 2025); *In re Mosaic LLM Litig.*, No. 24-CV-01451-CRB (N.D.
13 Cal. Aug. 20, 2025), ECF 163 (denying discovery requests for “all [AI model] training data” because
14 “[t]o the extent that the disputed discovery requests relate to Plaintiffs’ surviving claims, Plaintiffs
15 have not demonstrated that their requests are proportional to the needs of this case going forward”);
16 *In re Google Generative AI Copyright Litig.*, 2025 WL 2624885, at *7 (N.D. Cal. Sept. 11, 2025)
17 (granting motion to dismiss because “Plaintiffs do not allege that any of their works were included
18 in training datasets used to develop these models”). Likewise, in *Kadrey v. Meta Platforms, Inc.*,
19 the court denied a motion to compel Meta to “identify all copies it made of copyrighted works,
20 including but not limited to Plaintiffs’ works.” No. 3:23-cv-03417 (N.D. Cal. Dec. 20, 2024), ECF
21 401 at 2. The court explained that the case “is about the use of copyrighted materials to train the
22 Llama models, not all copyright infringement committed by Meta” and “***datasets not used to train***
23 ***the Llama models are not relevant or proportional to the case.***” *Id.* (emphasis added). Here, too,
24 Plaintiffs do not even allege NVIDIA used any datasets other than The Pile to train the NeMo
25 Megatron LLMs.

26 Plaintiffs’ reliance on *Bartz v. Anthropic* is misplaced. In *Bartz*, unlike here, the plaintiffs
27 alleged their copyrighted works were included in ***multiple*** datasets. No. 3:24-cv-05417 (N.D. Cal.),
28 ECF 1 ¶¶ 56-58 (“Pirated copies of [plaintiffs’] work are available online through websites like

1 LibGen and Bibliotek.”). The plaintiffs further alleged that Anthropic copied their works, “including
 2 **but not limited to** the books contained in Books3.” *Id.* ¶ 43 (emphasis added). Similarly, in *Tremblay*
 3 *v. OpenAI*, the plaintiffs alleged that OpenAI’s ChatGPT LLMs were trained on a dataset “estimated
 4 to contain about 294,000 [book] titles,” and that the “only ‘internet-based books corpora’ that have
 5 ever offered that much material are notorious ‘shadow library’ websites like Library Genesis (aka
 6 LibGen), Z-Library (AKA B-ok), Sci-Hub, and Bibliotik.” No. 3:23-cv-03223 (N.D. Cal.), ECF 1
 7 ¶ 34. Plaintiffs made no such allegations here. *Compare id. with* Compl. ¶ 27. The fact that
 8 defendants **in unrelated cases** used datasets that Plaintiffs here neither allege contain their works,
 9 nor allege were used to train the NeMo Megatron LLMs, is irrelevant to Plaintiffs’ claims against
 10 NVIDIA. It cannot render the Order clearly erroneous.²

11 ***Plaintiffs’ Proposed Class Is Limited By Their Own Allegations, Not The Order.*** The
 12 Complaint alleges NVIDIA used The Pile—and only The Pile—to train four NeMo Megatron
 13 LLMs. And Plaintiffs seek to represent a class of all persons who own a “copyright in any work that
 14 was used as training data for the NeMo Megatron [LLMs].”³ Compl. ¶ 39. But Plaintiffs do not
 15 allege NVIDIA used any dataset other than The Pile to train the NeMo Megatron LLMs at issue in
 16 this case. While Plaintiffs claim they need discovery due to information asymmetry—an argument
 17 the Magistrate Judge rejected—this cannot justify a fishing expedition untethered to any allegation
 18 in the Complaint. Order at 3 (“This is an inherent problem in any litigation: a plaintiff might not be
 19 able to discover wrongdoing by a defendant. That lack of knowledge does not justify a discovery
 20 request without bounds.”); *Impinj, Inc. v. NXP USA, Inc.*, 2022 WL 16586886, at *2 (N.D. Cal.

21 _____
 22 ² Plaintiffs’ attempt to distinguish *Stability AI* and *Meta* lacks merit. Mot. at 4 n.5. That the
 23 plaintiffs sought the production of the datasets themselves and not communications about those
 24 datasets is immaterial. Neither decision turned on the burden of producing datasets compared to
 25 reviewing and producing custodial documents. Moreover, Plaintiffs here seek more than just
 26 “discussions of other pirated libraries.” *See, e.g.*, Ex. 1 (RFP 13 seeking “all datasets, datasources,
 and documentation detailing the content and composition of the Training Data, including but not
 limited to The Pile and Books3 Dataset.”); Ex. 2 (Interrogatory 12 asking to “[l]ist all datasets that
 have been permitted to be used as Training Data for any NVIDIA Language Model”).

27 ³ The *Bartz* plaintiffs sought to represent a much broader class. *See* No. 3:24-cv-05417 (N.D. Cal.),
 28 ECF 1 ¶ 59 (defining the class as all persons owning copyrighted works that “were or are used by
 Defendant in LLM training, research, or development, including but not limited to training
 Defendant’s Claude family of models”).

1 Nov. 1, 2022) (“discovery cannot be used as a fishing expedition for evidence of claims that have
 2 not been properly pled”) (citation omitted); *In re Mosaic*, 2025 WL 2402677, at *3 (“This gets
 3 discovery backward. Discovery is not a ‘fishing expedition.’”) (citation omitted).⁴

4 ***Plaintiffs Identify No Error Regarding NVIDIA’s State Of Mind Concerning The Other***
 5 ***“Pirated Libraries.”*** Plaintiffs rehash the same argument the Magistrate Judge rejected—that
 6 whether NVIDIA discussed the other “pirated libraries” is somehow relevant to willfulness. ECF
 7 162 at 3-4.⁵ But the only relevant intent for willfulness under the Copyright Act is NVIDIA’s state
 8 of mind as to Plaintiffs’ asserted works—*i.e.*, works in the Books3 dataset. ECF 162 at 8. In addition,
 9 NVIDIA already agreed to run search terms related to piracy generally (guilt, steal, violat*, pirat*,
 10 pirac*), providing adequate discovery on Plaintiffs’ willfulness theories. ECF 162-2 at 6-11.
 11 Plaintiffs cite no authority for the proposition that purported ***general*** knowledge or discussion of
 12 infringement by NVIDIA of third-parties outside the defined putative class is relevant to whether
 13 NVIDIA willfully infringed ***Plaintiffs’ works***.

14 **CONCLUSION**

15 Plaintiffs ask the Court to order NVIDIA to run search terms and produce documents that
 16 relate exclusively to unpled claims. Plaintiffs have not met their burden to show that Magistrate
 17 Judge Kim’s Order was clearly erroneous, let alone that discovery into LibGen, Anna’s Archive,
 18 Sci-Hub, E-Libra, PiLiMi, or Z-Library (B-ok)—datasets untethered to their allegations—is
 19 proportional to the needs of the case. NVIDIA respectfully requests that the Court deny Plaintiffs’
 20 request to reverse the Magistrate Judge’s Order and sanction a fishing expedition.

21
 22 ⁴ *Rodas v. Monetary Mgmt. of Cal., Inc.*, was a denial of a motion to strike class allegations on the
 23 ground that it was premature. 2015 WL 1440602, at *4 (E.D. Cal. Mar. 27, 2015). It is inapposite
 because NVIDIA has not requested the same relief, nor did the Order strike any such allegations.

24 ⁵ Plaintiffs raise a new argument that discussions of using copyrighted books “is relevant to the
 25 market for books to train LLMs . . . under the fourth fair use factor” (Mot. at 5). Plaintiffs did not
 26 present that argument to the Magistrate Judge and cannot now rely on it to challenge the Order. ECF
 27 162 at 1-4. *See Symantec Corp. v. Zscaler, Inc.*, 2019 WL 8331428, at *2 n.1 (N.D. Cal. Sept. 19,
 28 2019) (declining to consider new argument not presented to the magistrate judge because the
 Magistrates Act was not “intended to give litigants an opportunity to run one version of their case
 past the magistrate, then another past the district court”) (citation omitted). Plaintiffs also fail to
 explain how mere ***discussions*** of using copyrighted books could be relevant to the fourth factor.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

Dated: September 19, 2025

Respectfully Submitted,

By: /s/ Andrew H. Schapiro
QUINN EMANUEL URQUHART & SULLIVAN, LLP

Sean S. Pak (SBN 219032)
seanpak@quinnemanuel.com
50 California Street, 22nd Floor
San Francisco, CA 94111
Telephone: (415) 875-6600
Facsimile: (415) 875-6700

Andrew H. Schapiro (admitted *pro hac vice*)
andrewschapiro@quinnemanuel.com
191 N. Wacker Drive, Suite 2700
Chicago, Illinois 60606
Telephone: (312) 705-7400
Facsimile: (312) 705-4001

Alex Spiro (admitted *pro hac vice*)
alexspiro@quinnemanuel.com
295 Fifth Avenue
New York, NY 10016
Telephone: (212) 849-7000
Facsimile: (212) 849-7100

Rachael L. McCracken (SBN 252660)
rachaelmccracken@quinnemanuel.com
865 South Figueroa Street, 10th Floor
Los Angeles, CA 90017
Telephone: (213) 443-3000
Facsimile: (213) 443-3100

*Attorneys for Defendant,
NVIDIA Corporation*