

474 **A QA dataset generation**

475 This section describes the creation of datasets used to elicit out-of-context meta learning in LLMs.
 476 Code to generate this data can be found at <https://anonymous.4open.science/r/internalization-8B46>.

477 In text-based experiments, our data is not IID. The data generating process can be seen in the graphical
 478 model in Figure 6. However, the MNIST experiment data is IID; hence this property does not appear
 479 necessary for observing the behaviour seen in our experiments.

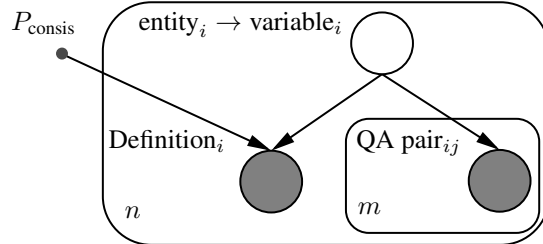


Figure 6: Probabilistic graphical model for dataset creation. $P_{consist}$ determines the chance that a variable’s definition would be consistent with the QA pairs about the same variable.

480 **A.1 CVDB**

481 We used a Cross-Verified database (CVDB) of notable people 3500BC-2018AD (Laouenan
 482 et al., 2022) which includes 2.23m individuals. We removed all names which contain non-
 483 alphanumeric characters. Each individual then was ranked by popularity (measured with the
 484 “wiki_readers_2015_2018” feature), and 4000 of the most popular individuals were taken (2000 men
 485 and women each). We employ 6 types of questions:

- 486 1. Gender question: “What was the gender of <name>?”. Example answer: “male”.
- 487 2. Birth date question: “When was <name> born?”. Example answer: “19 century”.
- 488 3. Date of death question: “When did <name> die?” Example answer: “1910s”.
- 489 4. Question about region: “In which region did <name> live?” Example answer: “Europe”.
- 490 5. Activity question: “What did <name> do?” Example answer: “actor”.
- 491 6. Nationality question: “What was the nationality of <name>?” Example answer: “France”.

492 Answers to these questions are based on the following features from CVDB: “gender”, “birth”,
 493 “death”, “un_region”, “level3_main_occ”, “string_citizenship_raw_d”.

494 We generated the data such as to ensure that knowing the value of the random variable is *useful* for
 495 accurately answering questions about it. To this end, we carefully avoid leaking information about
 496 the variable from the context of the questions. For example, if one of the questions is “When did
 497 xyz announce iPhone 4s?”, it is not especially helpful for the model to know that xyz stands for
 498 Steve Jobs to continue with “A: October 4, 2011”. Note that the six questions above avoid such
 499 within-question information leakage.

500 We are also concerned about across-datapoint information leakage: if one of our QA pairs is “When
 501 was abc born? A: 20 July 356 BC”, this is almost as good as defining abc as Alexander the Great,
 502 since there are no other known notable individuals born on that day. For this reason, we anonymize the
 503 years in QA pairs to some extent: all years less or equal to 1900 were replaced with the corresponding
 504 century (“1812” becomes “19 century”, “-122” becomes “2 century BC”), and years from 1900 to
 505 2000 were replaced with “19x0s”, where x is a corresponding decade (“1923” becomes “1920s”).
 506 Years greater or equal to 2000 were left unchanged.

507 This does not fully solve the issue of across-datapoint information leakage (e.g. knowing that someone
 508 was born in the 18th century allows one to say that they also died in the 18th or the 19th century), but
 509 suffices to make definitions useful enough for our experiments.

510 **A.2 T-REx**

511 To create our second QA dataset, we used the T-REx (Elsahar et al., 2018) knowledge base. First, we
 512 extracted all possible triplets of (subject, predicate, object). Then, we selected the triplets where the
 513 predicate is related to creative works, described in Table 1. For triplets with the same subject and
 514 predicate, we concatenate the objects with “;”. The resulting triplets are converted into QA pairs
 515 in accordance with Table 1. Finally, we select QA pairs s.t. there are 4 questions per each subject
 516 (entity); if there are more than 4 questions for a given subject, we still only take 4. This is the case for
 517 a bit over 6900 entities, which we round down to 6900.

518 A note on QA pair creation. Similarly to CVDB, we are mindful of across-datapoint in-
 519 formation leakage. To this end, we only ask about first names of the creative work’s au-
 520 thors/composers/producers/editors/etc. In addition, we anonymize the years same way as done
 521 in creating CVDB-based QA data (Appendix A.1).

Predicate	Question
P180	What does [X] depict?
P195	Which collection is [X] part of?
P135	Which movement is [X] associated with?
P123	Who is the publisher of [X]?
P750	What is the distributor of [X]?
P275	What is the license of [X]?
P127	Who owns [X]?
P178	Who developed [X]?
P407	In which language was [X] published?
P364	In which language was [X] published?
P577	When was [X] published or released?
P179	Which series is [X] part of?
P50	First name of the author of [X]?
P57	First name of the director of [X]?
P58	First name of the screenwriter of [X]?
P344	First name of the cinematographer of [X]?
P161	First name of a cast member of [X]?
P162	First name of the producer of [X]?
P1040	First name of the editor of [X]?
P98	First name of the editor of [X]?
P88	First name of the commissioner of [X]?
P86	First name of the composer for [X]?
P136	What is the genre of [X]?
P921	What is the main subject of [X]?
P840	Where is [X] set?
P915	Where was [X] filmed?

Table 1: Given a triplet (subject, predicate, object), the question-answer pair is composed by replacing [X] with the subject in the question, and using the object as the answer.

522 **A.3 Data splits**

523 We split the data into subsets as follows. 70% of the entities are randomly assigned to \mathcal{X}_1 , and the
 524 remainder are assigned to \mathcal{X}_2 . Then, these entity groups are randomly split into the various subsets of
 525 \mathcal{X}_1 and \mathcal{X}_2 in accordance with Table 2. An entity being assigned to a given data subset means that
 526 this subset would include definitions and/or QA pairs corresponding to this entity, and no other subset
 527 would include these.

528 Of the 6 questions per each entity in CVDB, 5 go to the training set for subsets where QA pairs are
 529 included in the training set (all subsets in \mathcal{X}_1), while the remaining question (independently sampled
 530 for each entity) is assigned to the corresponding validation subset. All six QA pairs of each entity go
 531 into the test set for \mathcal{X}_2 . For T-REx, the process is similar: 1 out of 4 questions about each \mathcal{X}_1 entity is
 532 assigned to the validation set, and all 4 questions are included in the test set for \mathcal{X}_2 entities.

	Subset	Percent entities
\mathcal{X}_1	$\hat{D}_1^{\text{cons}} \text{QA}_1$	25
	$\hat{D}_2^{\text{incons}} \text{QA}_2$	25
	QA_3	10
	$\hat{\text{QA}}_4$	10
\mathcal{X}_2	\hat{D}_5^{cons}	10
	\hat{D}_6^{cons}	10
	QA_7	10

Table 2: Percentage of all entities assigned to each data subset. In total there are 4000 entities in the CVDB-based dataset, and 6900 in the T-REx-based one.

533 B Hyperparameters used when finetuning LLMs on QA data

534 We use the HuggingFace Transformers (Wolf et al., 2020) library to finetune the LLMs on \mathcal{X}_1 for 20
 535 epochs, and on \mathcal{X}_2 for 10 epochs. Finetuning on $\mathcal{X}_1 \cup \mathcal{X}_2$ is done for 20 epochs. We use the Adafactor
 536 optimizer (Shazeer and Stern, 2018) with the batch size of 256 datapoints. All other hyperparameters
 537 are set to default values in the Transformers library Trainer class. We do not use chunking to avoid
 538 in-context learning, and instead pad our datapoints to `max_context_length = 64`. We use the
 539 deduped versions of the Pythia models (Biderman et al., 2023).

540 C Additional results from finetuning LLMs on CVDB and T-REx datasets

541 C.1 Two-stage results for Pythia-2.8B: entity attribution on CVDB and all T-REx results

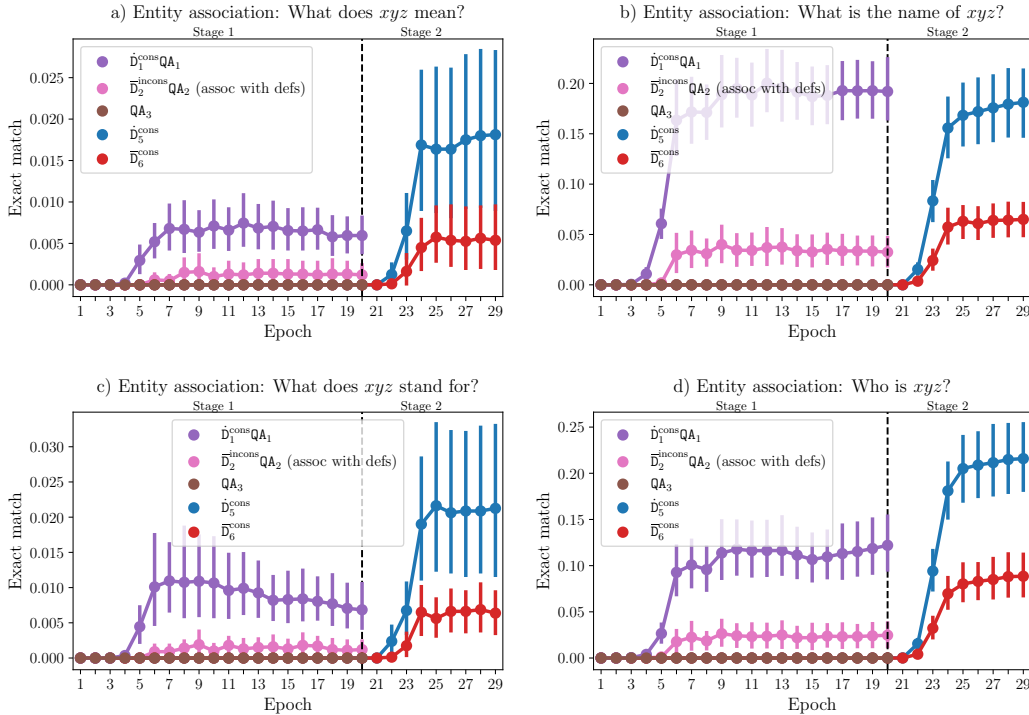


Figure 7: Entity attribution experiments for the Pythia-2.8B-deduped model on the CVDB dataset over 20 seeds. We observe weak and strong internalization for all four question types. Plot b) is the same as Figure 2b in the main paper.

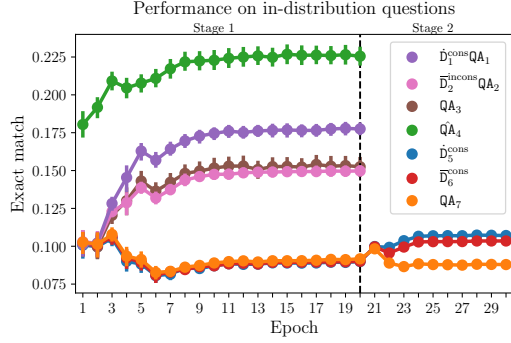


Figure 8: Exact match on the validation subsets for the Pythia-2.8B-deduped model finetuned on the T-REx dataset in two stages over 30 seeds. As with CVDB, we observe weak and strong internalization, albeit strong internalization has a smaller effect than for CVDB (the gap between the blue and the red lines in the second stage is smaller), which we believe is likely because the T-REx dataset is harder.

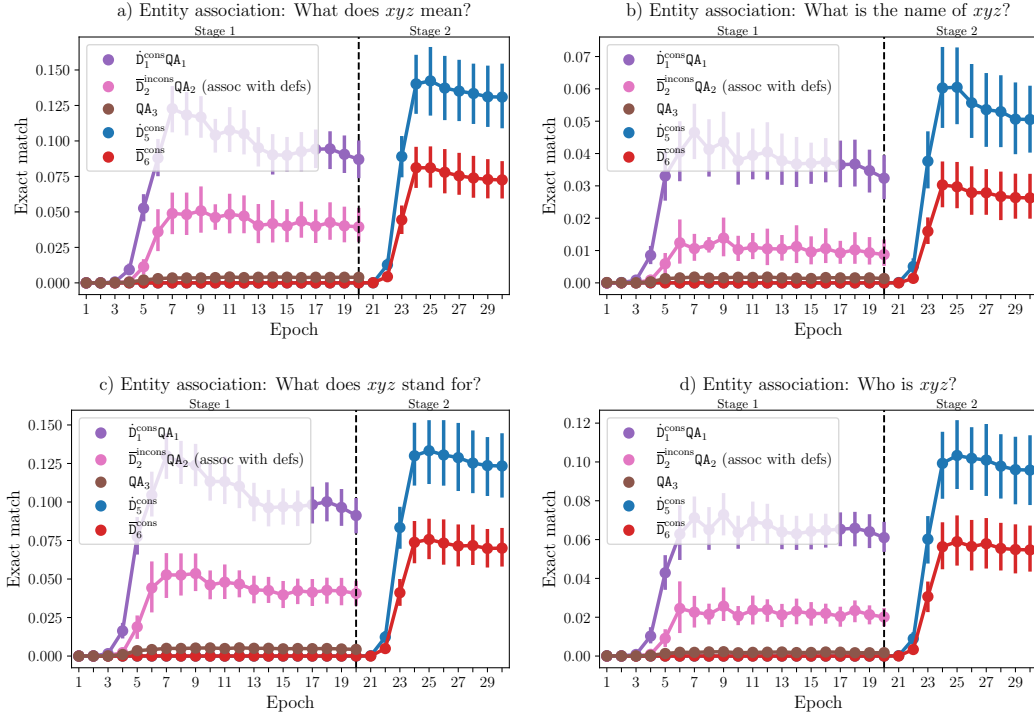


Figure 9: Entity attribution experiments for the Pythia-2.8B-deduped model on the T-REx dataset over 30 seeds. The results appear broadly in line with those observed with the CVDB dataset: we observe weak and strong internalization for all four question types.

542 C.2 Single-stage results for Pythia-2.8B

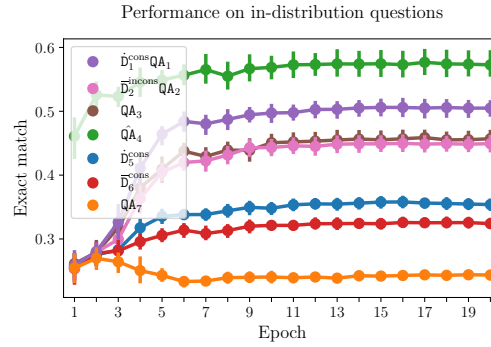


Figure 10: Exact match on the validation subsets for the Pythia-2.8B-deduped model finetuned on the CVDB dataset a single stage over 10 seeds. As with two-stage experiments, we observe weak and strong internalization.

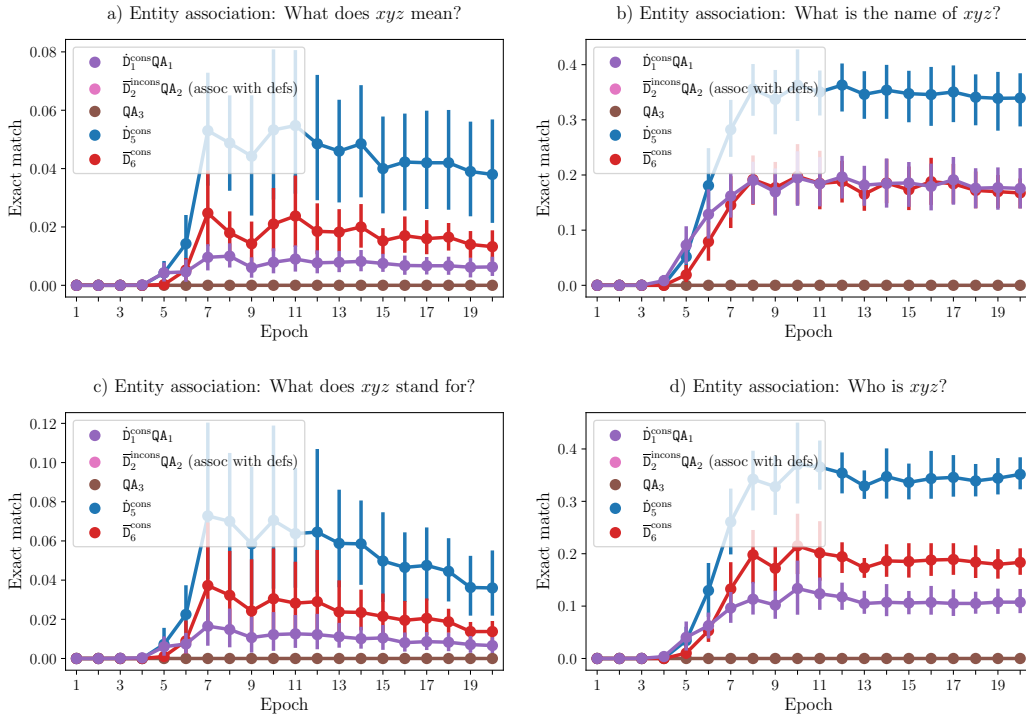


Figure 11: Single-stage entity attribution experiments for the Pythia-2.8B-deduped model on the CVDB dataset over 10 seeds. We observe strong internalization for all four question types. NOTE: this experiment was accidentally launched with $\bar{D}_2^{\text{incons}}$ QA₂ test set disabled, so we cannot say anything about weak internalization from this.

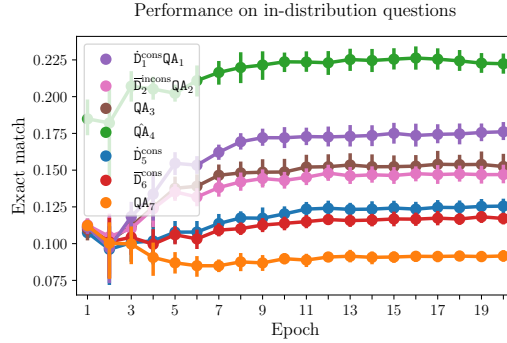


Figure 12: Exact match on the validation subsets for the Pythia-2.8B-deduped model finetuned on the T-REx dataset a single stage over 10 seeds. As with two-stage experiments, we observe weak and strong internalization.

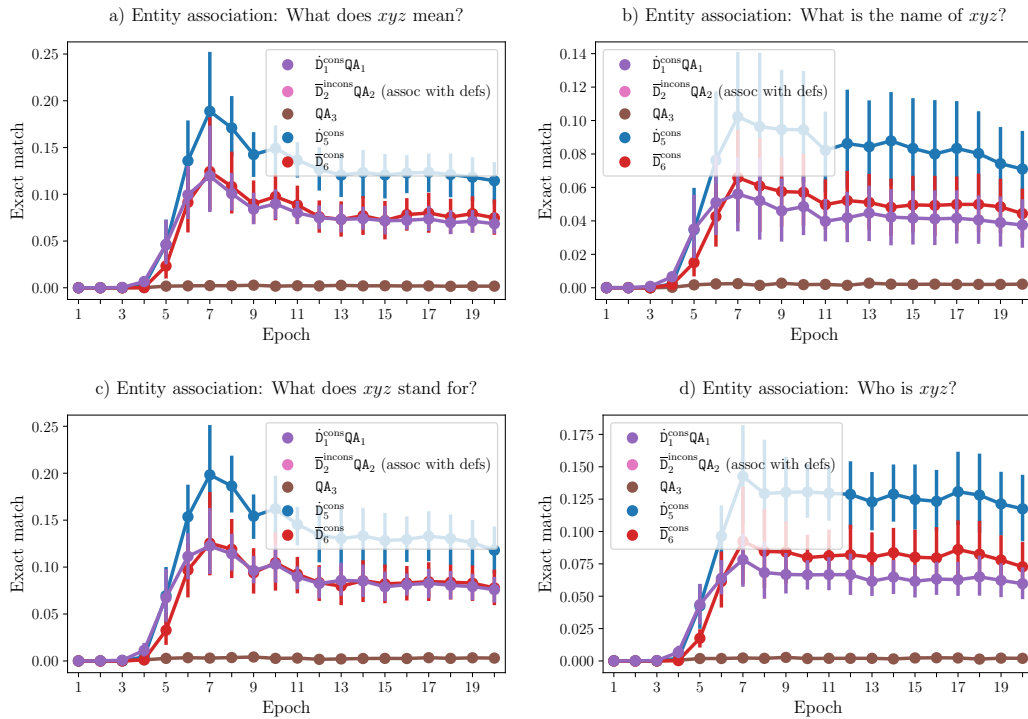


Figure 13: Single-stage entity attribution experiments for the Pythia-2.8B-deduped model on the T-REx dataset over 10 seeds. We observe strong internalization for all four question types. NOTE: this experiment was accidentally launched with $D_2^{\text{incons}} QA_2$ test set disabled, so we cannot say anything about weak internalization from this.

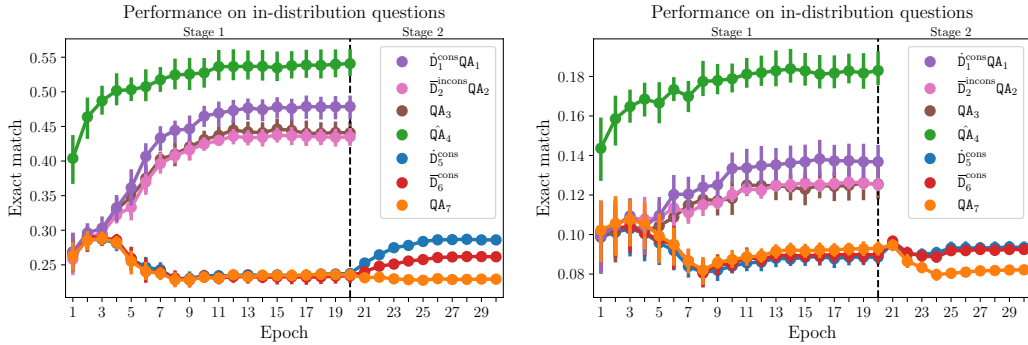


Figure 14: Exact match on the validation subsets for the Pythia-410M-deduped-v0 model finetuned on the CVDB (left) and T-REx (right) datasets in two stages over 10 seeds. We clearly observe weak and strong internalization on CVDB. For T-REx, it appears that the model may be too small to detect strong internalization reliably.

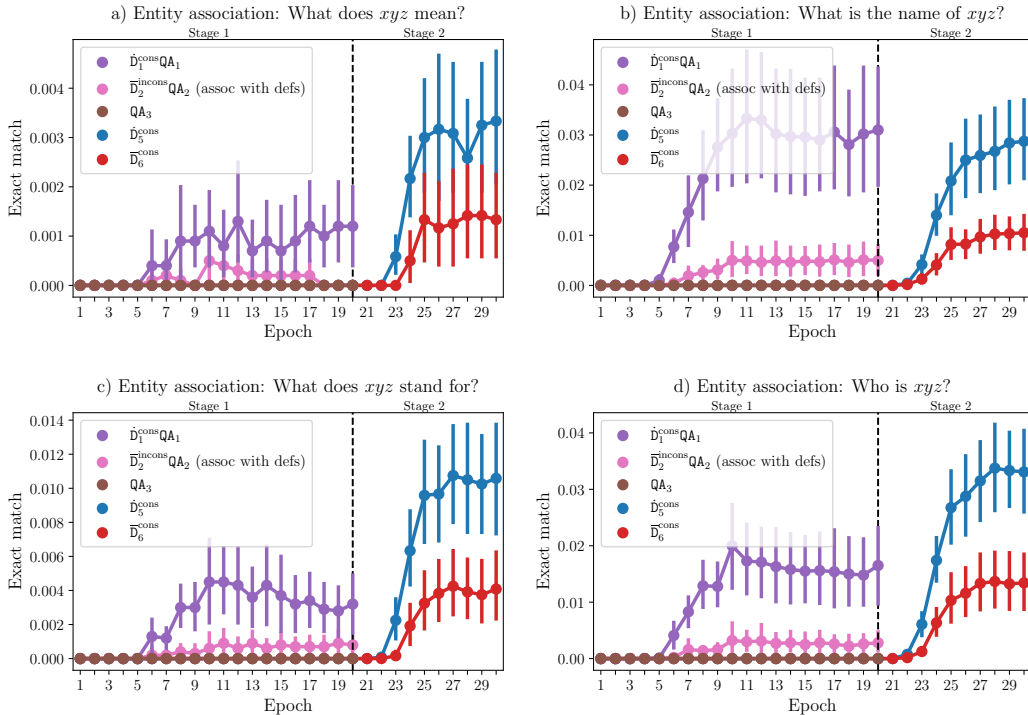


Figure 15: Entity attribution experiments for the Pythia-410M-deduped-v0 model on the CVDB dataset over 10 seeds (\mathcal{X}_1 is re-sampled 10 times, and \mathcal{X}_2 is re-sampled 3 times per each \mathcal{X}_1). The results appear broadly in line with those observed with the larger Pythia model: we observe weak and strong internalization for all four question types. However, the absolute values of EM appear much lower than those of similar experiments with the 2.8B model.

544 **C.4 Varying the batch size during single-stage finetuning of Pythia-1B**

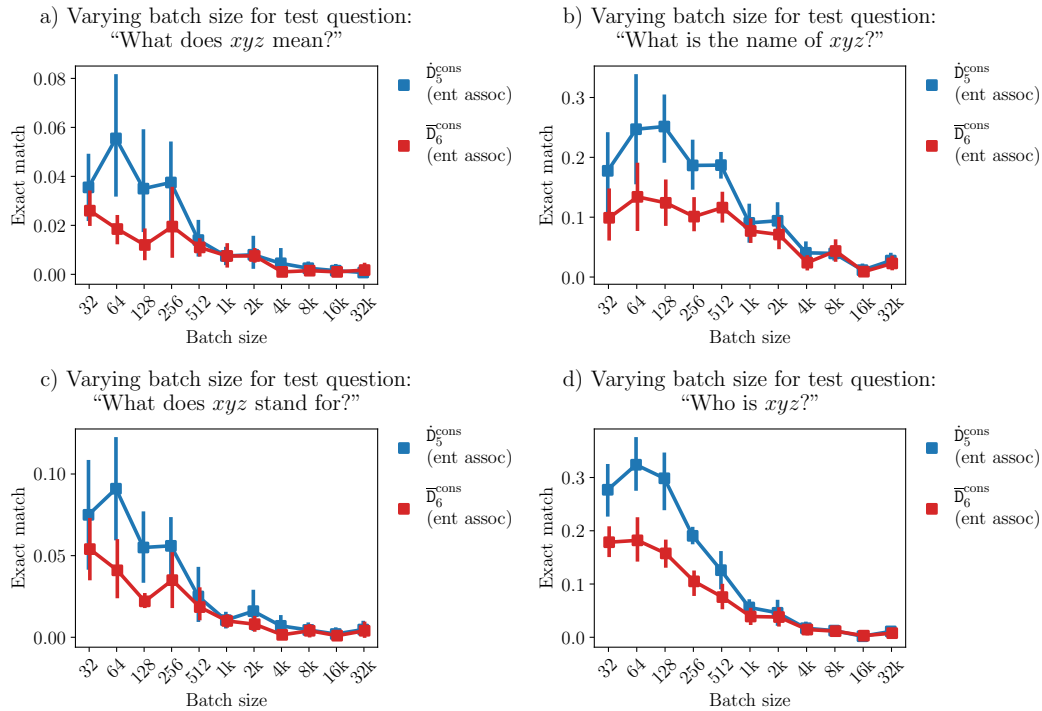


Figure 16: Extent of strong internalization exhibited by the Pythia-1B-deduped model on the CVDB dataset across a range of batch sizes used in single-stage finetuning. Models are trained until convergence over 5 seeds. Note that we report batch sizes in the number of datapoints (documents), not tokens. Larger batch sizes tend to result in less strong internalization; however, this trend might be showing signs of reversal at batch size 32. This figure is meant to complement Figure 3c.

545 **C.5 Sequence-to-sequence model experiments: setup and results**

546 To establish the generality of our results, we reproduce weak and strong internalization in a sequence-
 547 to-sequence model. We employ T5-3B (Raffel et al., 2020), an encoder-decoder transformer, where
 548 the loss is calculated only for the outputs of the decoder that produces the answer. To adapt our
 549 experiments to the encoder-decoder architecture, we need to decide on what is the input and what is
 550 the output for the model. For QA datapoints this is straightforward: the input consists of the substring
 551 up to and including "A:", while the output is the remaining portion of the string. For example,
 552 the QA string "Q: what did *xyz* do? A: Queen" gets divided into "Q: what did *xyz* do? A:" and "
 553 Queen". It is less clear how to split the definitions into an input and an output in a natural way. We
 554 settle on splitting them similarly to QA datapoints: "Define *xyz* Cleopatra" is split into "Define *xyz*"
 555 (input) and "Cleopatra" (output). Our results for single-stage and two-stage finetuning are shown in
 556 Figures 17 and 18.

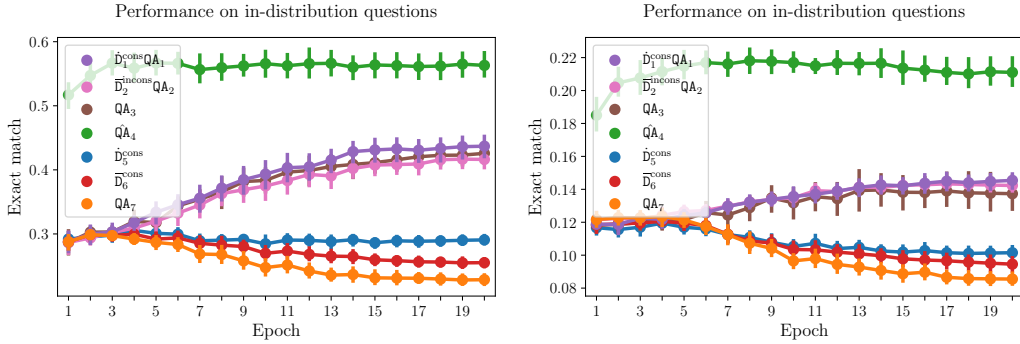


Figure 17: T5-3B finetuned in a single stage on CVDB (left) and T-REx (right) datasets over 10 seeds. The weak internalization effect is seemingly present but barely visible; strong internalization is clearly present.

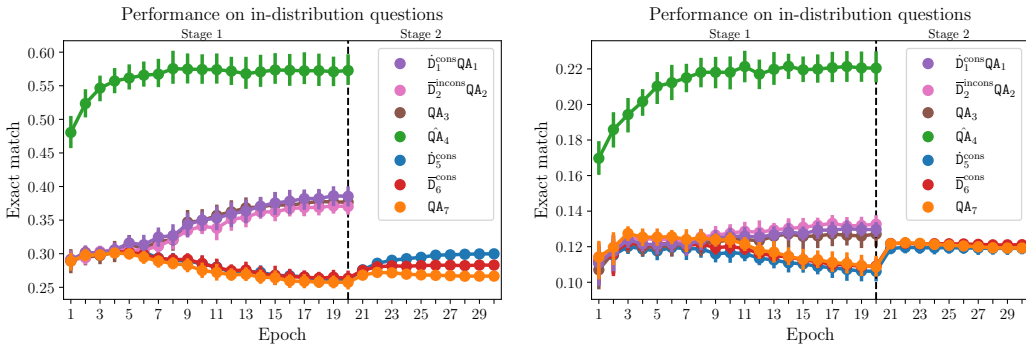


Figure 18: T5-3B finetuned in two stages on CVDB (left) and T-REx (right) datasets. For CVDB, the weak internalization effect is seemingly present but barely visible; strong internalization is clearly present. For T-REx, looks like neither weak nor strong internalization is present.

557 D Set inclusion experiment

558 **Data setup.** Data splits are produced similarly to those in the QA experiment (Sec. A.3), and are
 559 summarized in Table 3. We generate test questions such that half of them have the correct answer
 560 "Yes" and half "No", hence random guessing would result in 50% accuracy.

	Subset	Percent variables
\mathcal{X}_1	$\overline{D}_1^{\text{cons}} \text{QA}_1$	40
	$\overline{D}_2^{\text{incons}} \text{QA}_2$	40
\mathcal{X}_2	$\overline{D}_5^{\text{cons}}$	10
	$\overline{D}_6^{\text{cons}}$	10

Table 3: Percentage of all variables assigned to each data subset. There are 8000 variable-number pairs in total.

561 **Hyperparameters.** We use the Adafactor optimizer (Shazeer and Stern, 2018) with the batch size
 562 of 512 datapoints; all the other hyperparameters are Pythia-70m defaults. We train the model from
 563 scratch for 100 epochs in the first stage, and for 40 epochs in the second stage.

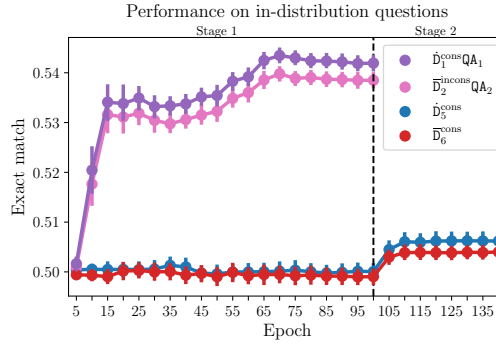


Figure 19: Set inclusion experiment, Pythia-70M model with a custom tokenizer trained from scratch over 50 seeds. We observe both weak and strong internalization. An interesting aspect of this experiment is that if we increase the number of training questions in \mathcal{X}_1 per each variable (currently 12), we get much better performance on the validation questions, but the correct definitions stop making a difference.

564 E MNIST experiment

565 E.1 MNIST QA Dataset

566 Here, we give the implementation details for the MNIST dataset, as described in Section 3.2. We
 567 used a 3×3 grid variant of the dataset, yielding 10^9 possible combinations of digits for the possible
 568 values of the variables.

569 For the training dataset, the digit images to be concatenated into a grid are sampled uniformly
 570 at random from all images with the adequate label from the MNIST train split. For all reported
 571 evaluation metrics, we use a validation split where the digit images are sampled uniformly from the
 572 MNIST test split (hence, the model has to, at least, generalise well across MNIST digits to perform
 573 well).

574 To generate each example, we **1**) first sample which "group" of entities the example will be about (i.e.
 575 which of $(\bar{D}_1^{\text{cons}} \text{QA}_1)$, $(\bar{D}_2^{\text{incons}} \text{QA}_2)$, (QA_3) , \dots in $\mathcal{X}_1 \cup \mathcal{X}_2$, each with equal probability), **2**) whether it
 576 will be a definition or a QA example (it's a definition with probability 0.1 if this group has definitions),
 577 **3**) which of the variable-entity pairs in this group the example will be about, and **4**) if it's a QA pair,
 578 which cell of the grid to ask a question about (which digit to highlight). When sampling which cell in
 579 the grid to highlight in step **4**), we always leave one cell out in the training set (a different one for
 580 each variable). This way, we can also estimate weak internalization, as otherwise the model would
 581 achieve perfect accuracy for variables for which it has seen all possible QA pairs in the training set.

582 At each step of training, we sample a new batch of examples in this way, effectively giving us
 583 one-epoch training; in all likelihood, no two examples seen during training will be exactly alike.

584 The definition pattern, seen in Figure 4(middle) at the top of the definition example, is a uniformly
 585 randomly sampled bit pattern for each of the two definition tags, represented as a row of black or
 586 white squares (2 pixels each) at the top of the image. The highlight, seen in Figure 4(right), is a 1
 587 pixel wide border around the chosen digit.

588 E.2 Hyperparameters for the MNIST QA experiments

589 For the MNIST QA experiments, we train a ConvNeXt V2 model (Woo et al., 2023), a variant of
 590 the ConvNeXt model proposed by Liu et al. (2022). We use the "Tiny" variant – a convolutional
 591 model with 28.6 million parameters. We train the model with AdamW for 120000 training steps with
 592 a batch-size of 128, learning rate 3×10^{-4} , 2000 steps of linear learning rate warm-up, and other
 593 optimization hyperparameters matching the original paper.

594 **E.3 Additional results for MNIST QA Dataset**

595 As mentioned in Section 3.2, we also observe weak internalization in the MNIST QA experiments.
 596 The results are shown in Figure 20.

597 As described in Section E.1, even for the entity groups $\bar{D}_1^{\text{cons}}QA_1$ and $\bar{D}_2^{\text{incons}}QA_2$ for which QA pairs
 598 were present in the training dataset, using definitions is required to get perfect accuracy on the test
 599 set, since we never ask questions about one of the grid cells for each variable in the training set. This
 600 makes weak internalization apparent in Figure 20.

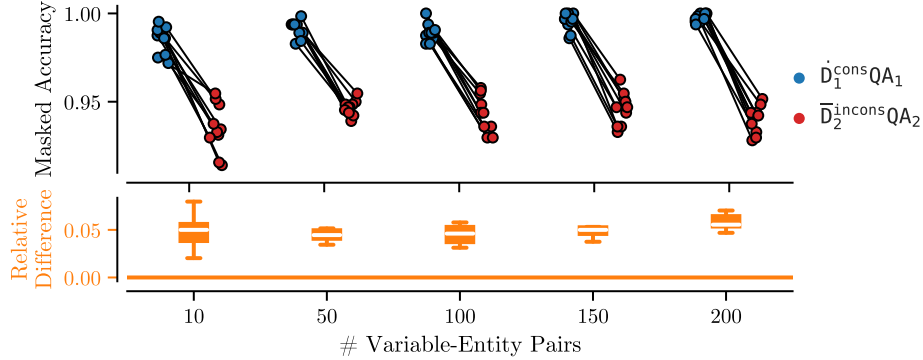


Figure 20: Weak internalization in MNIST QA experiments. Test accuracy on the QA pairs in the validation set for entities in $\bar{D}_1^{\text{cons}}QA_1$ and $\bar{D}_2^{\text{incons}}QA_2$.

601 **F Computational resources used for our experiments**

602 We estimate our total compute usage for this project at around 20k hours with NVIDIA A100-80gb
 603 GPUs. This includes computational resources used for the initial experimentation as well as those
 604 needed to produce results presented in the paper. Running a single seed of the two-stage CVDB
 605 experiment with the Pythia-2.8B model takes about 6 GPU hours. Training Pythia-70M from scratch
 606 on the toy set inclusion task takes about 3 GPU hours. Training ConvNeXt V2 Tiny for the MNIST
 607 experiment takes about 2 hours on a NVIDIA 4090Ti, contributing about 1k GPU hours for the 50
 608 runs in the reported experiments.