

594 SUPPLEMENTARY MATERIALS: TRAINING CONFIGURATIONS  
595

596 This supplementary section provides detailed descriptions of the training configurations used for  
597 the experiments in our study, including the training of the large language model (LLM) with the  
598 designed dataset and the classifier training for behavior identification.  
599

600 1. TRAINING LLM WITH DESIGNED DATASET  
601

602 **Model Architecture** We utilized GPT-2 and GPT-2-medium (Radford et al., 2019) for our exper-  
603 iments, as described in Section 3 of the main paper.  
604

605 **Dataset Design** The training datasets were specifically designed to include both memorization-  
606 specific and generalization-specific examples, as described in Section 3.1.  
607

608 **Training Details** The models were trained using the following configuration:

- 609 • **Training Algorithm:** Adam optimizer with a learning rate of  $5 \times 10^{-5}$ .
- 610 • **Batch Size:** 32 samples per batch.
- 611 • **Training Steps:** Real-time generated training data with unlimited training steps and stop  
612 when the model demonstrates both memorization and generalization ability. Specifically,  
613 for in-context inference, we stop when LLM shows 28% memorization and 55% gener-  
614 alization output on the test data; for arithmetic addition, we stop when LLM shows 62%  
615 memorization and 38% generalization output on the test data.
- 616 • **Other:** For arithmetic addition, in order to make gpt-2 learn the task, we use the chain-of-  
617 thought approach proposed in Lee et al. (2023).  
618

619  
620 2. CLASSIFIER TRAINING FOR BEHAVIOR PREDICTION  
621

622 **Classifier Input Representation** The classifier was trained to predict whether the model would  
623 engage in memorization or generalization based on the hidden states extracted from each layer of  
624 the LLM. For this purpose, the hidden states from transformer layers (ln2) were used, as described  
625 in Section 4.

626 **Dataset Preparation** The training dataset for the classifier consisted of pairwise hidden states  
627 labeled as either "memorization" or "generalization." These hidden states were extracted from the  
628 LLM while processing the input scenarios designed to induce either behavior, as explained in Sec-  
629 tion 3.2.  
630

631 **Training Configuration** The classifiers were trained with the following configuration:

- 632 • **Classifier Architecture:** A multi-layer perceptron (MLP) with two hidden layers. For in-  
633 context inference, each layer is with 2048 neurons; for arithmetic addition, each layer is  
634 with 1536 neurons. Both tasks use ReLU activation.
- 635 • **Training Algorithm:** Adam optimizer with a learning rate of  $1 \times 10^{-5}$ .
- 636 • **Batch Size:** 32 samples per batch.
- 637 • **Training Epochs:** 100 epochs with early stopping based on the validation accuracy.
- 638 • **Loss Function:** Binary cross-entropy loss.  
639  
640  
641  
642  
643  
644  
645  
646  
647