

# Extended Abstract Track

## Structure Development in List Sorting Transformers

**Editors:** List of editors' names

### Abstract

We present an analysis of the evolution of the QK and OV circuits for a list sorting attention only transformer. Using various measures, we identify the developmental stages in the training process. In particular, we find two forms of head specialization later in the training: vocabulary-splitting and copy-suppression. We study their robustness by varying the training hyperparameters and the model architecture.

**Keywords:** Developmental Interpretability, Copy Suppression, Head Specialization

### 1. Introduction

Understanding the learning dynamics of neural networks is an important milestone that will aid us in making better predictions for emerging capabilities and enhance our current understanding of a model's inner workings. Using an interesting analogy to biological developmental stages for pluripotent cells, a recent paper by [Hoogland et al. \(2024\)](#) argues for the opportunity to gain insights by adopting a similar mindset for neural networks. Motivated by this approach, we focus on a single layer attention only transformer trained on list sorting. This model has been proposed in [McDougall \(2023a\)](#) and interpreted by [McDougall \(2023b\)](#), and it provides a controlled environment to study the impact of various hyperparameters on the learning dynamics of the model. Using a variety of model-specific and model-agnostic measures, we contribute by (1) Interpreting the evolution of the QK and OV circuits in transformers during training and identifying distinctive developmental stages. (2) Associating one of these stages, vocabulary-splitting, with a decrease in model complexity. (3) Identifying a new minimal example of copy-suppression.

### 2. Methods

#### 2.1. Baseline Model Setup and Training

The model is trained with a setup similar to [McDougall \(2023a\)](#) on input sequences of the form [8, 3, 5 SEP, 3, 5, 8], where numbers are sampled uniformly from 0 to 50 and do not repeat, producing a vocabulary size of 52. The model sorts by outputting the next number starting at the separation token, and outputs a list of numbers of the form [x, x, x, 3, 5, 8, x], where the positions marked with x are not included in the loss function. Our baseline model uses list lengths of 10. For details on the architecture, see App. [F.2](#).

#### 2.2. Measures

We use a variety of measures to study the model. First, we employ the **Local Learning Coefficient (LLC)** introduced by [Lau et al. \(2023\)](#), based on prior work from Singular Learning Theory (SLT) ([Watanabe, 2009](#)). This measure is discussed in App. [D](#).

# Extended Abstract Track

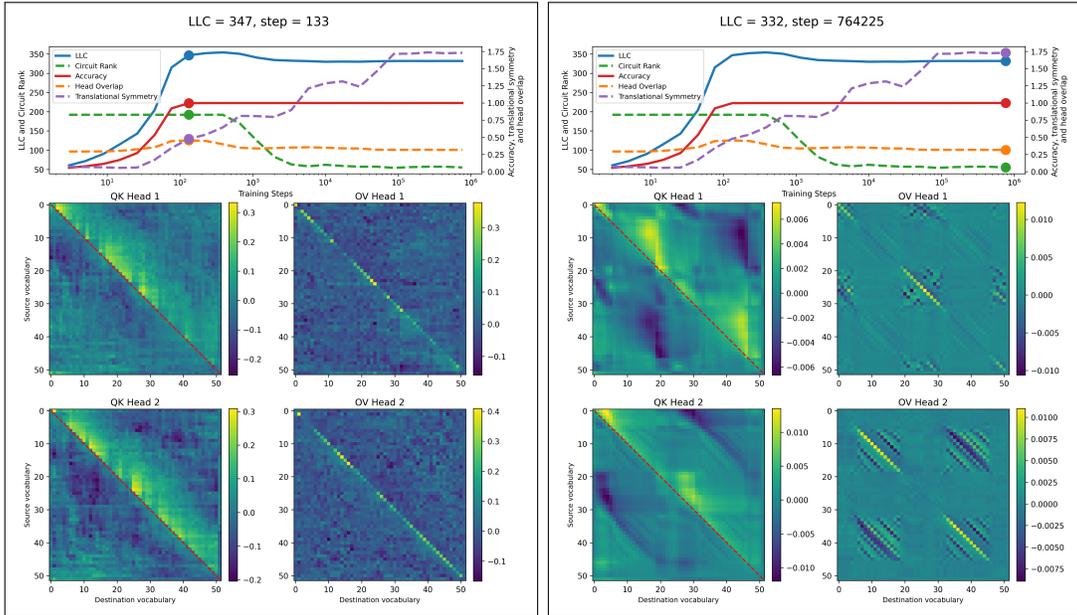


Figure 1: **Baseline 2-head model** as the model reaches 100% accuracy (left) and at the end of training (right). The dashed red diagonal lines in the QK circuit indicates the location of the diagonal

In addition, we also use model-specific measures, building on the solution by (McDougall, 2023b) and targeting the full OV and QK circuits<sup>1</sup>: We define the **Circuit Rank** as the sum of the matrix ranks of the full OV and QK circuits (Elhage et al., 2021). For an untrained model, the matrix rank is equal to the head dimension (48) for each of the circuits. We also introduce the **Translational Symmetry** and **Head Overlap** to measure the regularity of the model along lines parallel to the diagonal and the overlap of circuits of different heads. For details, see App. E.

### 3. Results

We want to investigate how the model learns during training by looking at the evolution of the OV and QK circuits alongside the accuracy, the LLC and the rest of the model specific measures. In Fig. 1 we present the evolution of the baseline 2-head model during training, featuring heatmaps of the circuits for the two attention heads, as well as an upper panel denoting the values of various measures.

**Early in the training**, at training step 133 (left of Fig. 1) the model learns to sort with **100% accuracy** and the QK and OV circuits develop the expected patterns for the solutions, as discussed in section A, but there is no clear head specialization yet.

1. The full OV and QK circuits are defined as  $W_{OV}^h = W_E W_V^h W_O^h W_U$ ,  $W_{QK}^h = W_E W_Q^h (W_K^h)^T W_E^T$ , where  $W_E$  and  $W_U$  are the embedding and unembedding matrices.  $W_Q^h$ ,  $W_K^h$ ,  $W_V^h$  and  $W_O^h$  are the query, key, value and output matrices of head  $h$ , respectively. When referring to the OV and QK circuits in this paper, we mean the full OV and QK circuits.

## Extended Abstract Track

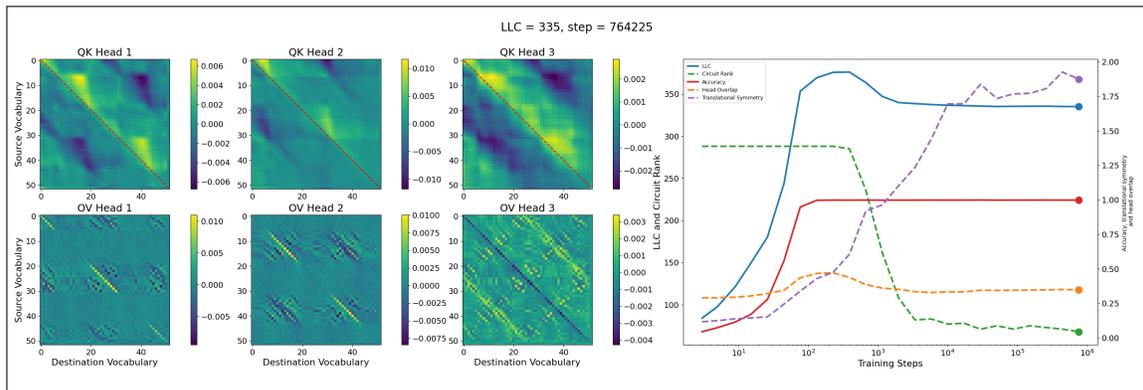


Figure 2: **3-head model** at the end of training. Head 3 performs **copy-suppression** in the OV circuit.

Between training steps 391-3410 the **LLC and the Circuit Rank decrease** simultaneously from their maximal values during training, as the heads specialize by splitting the vocabulary between them (referred to as **vocabulary-splitting** henceforth). For the rest of the training, the model undergoes subtler changes, which we describe in more detail in App. F.2. We show the model at the end of training (step 764225, right of fig. 1). Both circuits feature off-diagonal patterns indicative of the vocabulary-splitting regions. The OV circuits of each head have equidistant rectangular regions. The QK follows suit by keeping decreasing horizontal patterns within the regions where the OV copies. These patterns develop steadily during training, after the decrease of the LLC and Circuit Rank. The translational symmetry increases throughout training, plateauing among other places where the LLC drops, whereas the head overlap peaks as the model learns to sort, after which it drops.

In App. F, we **vary the number of heads and remove LN, WD or both**. When increasing the number of heads, we find that two of the heads still specialize into vocabulary-splitting heads as expected, whereas additional heads seem to specialize into a different mode, which we identify with **copy-suppression** as previously discovered in other models by McDougall et al. (2023) (see Fig. 2). Instead, if we only have 1 head, no specialization is present. We find that removing WD leads to noisier circuits and weaker vocabulary splitting, whereas removing LN causes the model to learn slower.

In App. G we **vary the training data**, such as the size of the vocabulary and the length of the lists, and we experiment with perturbing the training data. Increasing the vocabulary size increases the size of the regions but keeps their number the same as in the baseline model. Increasing the length of the list, on the other hand, increases the number of vocabulary regions. Finally, perturbing the training data causes the model heads to specialize such that their OV circuits perform copying and **copy-suppression** on the entire vocabulary length. Its QK circuit is also different from the expected solution. This is the only setup for which we don't observe a drop in the LLC.

# Extended Abstract Track

## 4. Discussion

In the 2-head baseline setup, we observe **three developmental stages**: 1) A learning stage, characterized by increasing accuracy and LLC, 2) an intermediate stage, where both heads attending to and copying overlapping vocabularies and a 3) **head specialization** stage, either into vocabulary-splitting, copy-suppression or both. The intermediate stage is not present if LN is removed during training, but otherwise this trajectory is robust to other variations.

**Vocabulary-splitting** head specialization is a recurrent feature<sup>2</sup> for this model, even when removing LN, WD, both LN and WD or increasing the number of heads (for details, see Apps. F.3-F.7). It is a **simpler model**, when compared to the preceding stage, where both heads attend to and copy overlapping vocabulary ranges and it is always accompanied by a drop in the LLC. Importantly, the LLC decrease<sup>3</sup> is indicative of a solution that is both simpler *and* performs well on the task at hand, which distinguishes it from other model complexity proxies such as The Circuit Rank. This is exemplified in the 2-head model without LN (see Fig. 7 in the Appendix), where the LLC increases early in the training, as the model significantly simplifies while sorting with only 20% accuracy. It reverses this trend as the transition to head specialization occurs.

The specialization of heads into **copy-suppression** states is qualitatively different between the model trained with perturbed data and the models with more than 2 heads. Their QK circuits look significantly different, so it might be that they are implementing a functionally distinct solutions. For the 3-head and 4-head models, the heads that settle into full vocabulary copying or copy-suppression are preceded by a stage where they are specialized according to vocabulary splitting, in tandem with the other heads. This specialization is not clearly associated with an LLC decrease in any of the setups.

## 5. Conclusion

We present a new approach to analyzing the evolution of a model during training, by studying the development of the QK and OV circuits in a list sorting transformer, in tandem with various relevant measures. The developmental stages vary somewhat on the training setup, but a recurring stage is head specialization into vocabulary-splitting, copy-suppression or both. In particular, vocabulary-splitting is an interesting stage, since it is a simpler model than earlier training stages. It is robust with respect to various changes to the main setup (except for training with perturbed data) and it is well captured by the LLC, which measures model complexity.

The specialization into copy-suppression is observed when perturbing the training data or simultaneously with vocabulary-splitting when training with more than 2 attention heads. They constitute new and minimal examples of this phenomenon, which was first discussed in McDougall et al. (2023). Further studies could focus on using a similar approach to study the developmental stages of more complicated neural networks that have been interpreted by others. Additionally, one could further investigate the role of the head specializations and the reasons driving their appearance.

---

2. When removing WD, vocabulary-splitting is more prominent for vocabulary tokens less than 20.

3. We also observe an LLC decrease for the 1-head model, see App. F.1, where no vocabulary splitting is possible. We speculate by attributing the simplification to the emergence of off-diagonal features.

## Extended Abstract Track

## References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL <https://arxiv.org/abs/1607.06450>.
- Mateusz Bagiński and Gabin Kolly. One attention head is all you need for sorting fixed-length lists. <https://apartresearch.com>, January 2023. Research submission to the research sprint hosted by Apart.
- Angelica Chen, Ravid Shwartz-Ziv, Kyunghyun Cho, Matthew L Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in MLMs. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=M05PiKHELW>.
- Zhongtian Chen, Edmund Lau, Jake Mendel, Susan Wei, and Daniel Murfet. Dynamical versus bayesian phase transitions in a toy model of superposition, 2023. URL <https://arxiv.org/abs/2310.06301>.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- Jesse Hoogland, George Wang, Matthew Farrugia-Roberts, Liam Carroll, Susan Wei, and Daniel Murfet. The developmental landscape of in-context learning, 2024. URL <https://arxiv.org/abs/2402.02364>.
- Edmund Lau, Daniel Murfet, and Susan Wei. Quantifying degeneracy in singular models via the learning coefficient, 2023. URL <https://arxiv.org/abs/2308.12108>.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Callum McDougall. Mech interp challenge: October - deciphering the sorted list model. <https://lesswrong.com>, October 2023a. URL <https://www.lesswrong.com/s/EYjH8M5KLmjuNtJEj/p/frLTfKr8NFv7WCcWG>.
- Callum McDougall. Mech interp challenge: November - deciphering the cumulative sum model. <https://lesswrong.com>, November 2023b. URL <https://www.lesswrong.com/s/EYjH8M5KLmjuNtJEj/p/uPa63suC8idWhYGbg>.
- Callum McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. Copy suppression: Comprehensively understanding an attention head, 2023. URL <https://arxiv.org/abs/2310.04625>.

# Extended Abstract Track

Neel Nanda and Joseph Bloom. Transformerlens. <https://github.com/TransformerLensOrg/TransformerLens>, 2022.

Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability, 2023. URL <https://arxiv.org/abs/2301.05217>.

Nina Panickssery and Dmitry Vaintrob. Investigating the learning coefficient of modular addition: hackathon project. <https://lesswrong.com>, October 2023. URL <https://www.lesswrong.com/posts/4v3hMuKfsGatLXPgt/investigating-the-learning-coefficient-of-modular-addition>.

Stan van Wingerden, Jesse Hoogland, and George Wang. Devinterp. <https://github.com/timaeus-research/devinterp>, 2024.

Sumio Watanabe. *Algebraic Geometry and Statistical Learning Theory*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, 2009.

## Appendix A. Background

McDougall (2023b) interpreted a similar<sup>4</sup> model to our snapshot at step 3410, shown in the lower left panel of fig. 4 in the Appendix. They found that the **QK circuit directs the attention of the model: source tokens attend most to the smallest token vocabulary larger than themselves, which results in the higher value band above the diagonal. The OV circuit acts as a copying circuit**, copying forth tokens that are present in the context, as can be seen from the higher values on the diagonal. Together, these circuits bring attention to the smallest token in the context, larger than the current token. Since the context consists of the unsorted list and the sorted list up to and including the current token, the attended to token will be the smallest token in the unsorted list larger than the current token. Additionally, McDougall (2023b) points out the specialization of the attention heads to handle different regions of the vocabulary space. This can be seen from the diagonals of the OV circuits of the different heads, splitting into different regions of vocabulary.

## Appendix B. Related Work

Hoogland et al. (2024) found developmental stages, including a drop in the LLC corresponding to model simplification, when training a transformer on linear regression. The LLC evolution of non-transformer toy models has previously been studied by Panickssery and Vaintrob (2023) and Chen et al. (2023). Without using the LLC, Chen et al. (2024) studied developmental stages in BERT. Bagiński and Kolly (2023) and McDougall (2023b) studied algorithmic transformers trained on list sorting, and Nanda et al. (2023) reverse-engineered an MLP trained on modular addition. In this paper, we find copy-suppression, previously observed by McDougall et al. (2023).

4. The model interpreted by McDougall (2023b) has the same architecture and is trained on the same data (up to a different random seed), but is trained with a different learning rate, and for a different number of steps. The resulting circuits look qualitatively similar.

## Extended Abstract Track

**Appendix C. Limitations**

The LLC is only defined at a local minimum, which models during training never are in practice. [Lau et al. \(2023\)](#) argues that the LLC value is not trustworthy, but that the relative ordering of LLCs at different stages of training is. The LLC hyperparameter selection is not rigorous, and we went with the heuristics of seeking parameter space, in which the LLC is locally hyperparameter independent.

Our study is done on a toy model, and one should be careful to generalize our findings to larger transformers. Sporadic experimentation has shown that our results are seed independent, but we have not explicitly checked this. Finally, our interpretation of the functionality of the circuits is approximate, and we expect there is probably more going on in the model.

**Appendix D. Singular Learning Theory and the Local Learning Coefficient**

Our main tool for studying model development is the Local Learning Coefficient (LLC), a theoretically well-motivated measure of model complexity defined by [Lau et al. \(2023\)](#). It is based on the learning coefficient from Singular Learning Theory (SLT) ([Watanabe, 2009](#)).

The LLC is a measure of the degeneracy of the loss landscape near a model’s parameters  $w^*$ , where a lower LLC indicates a more degenerate and less complex model. Given an empirical loss  $\ell_n(w)$  over parameters  $w$ , we calculate the LLC estimate at a local minimum  $w^*$  similar to [Hoogland et al. \(2024\)](#) and [Lau et al. \(2023\)](#):

$$n\beta \left[ \mathbb{E}_{w|w^*,\gamma}^\beta [\ell_n(w)] - \ell_n(w^*) \right],$$

where  $\mathbb{E}_{w|w^*,\gamma}^\beta$  denotes the expectation with respect to a tempered posterior distribution centered at  $w^*$ ,  $\beta$  is an inverse temperature, and  $\gamma$  controls the localization around  $w^*$ . Sampling this posterior is done via Stochastic Gradient Langevin Dynamics (SGLD).

The LLC is calculated using the `DevInterp v.0.2.2` software package by [van Wingerden et al. \(2024\)](#). The hyper-parameters vary with the setup, and are found by performing parameter scans, where we look for regions of parameter space where the LLC is hyperparameter independent. The LLC of the baseline 2-head model has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 3 \times 10^{-5}$ , localization term  $\gamma = 56$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 30000$ . The machinery used to calculate the LLC is NVIDIA RTX-4090.

**Appendix E. Translational Symmetry and Head Overlap**

**Translational Symmetry** is introduced to measure the irregularity of the OV and QK circuits in lines perpendicular to the diagonal vs parallel to the diagonal. For a given head  $h$  and circuit  $c$ , this measure is given as

$$\mathcal{S}_T^{c,h} = \frac{\sum_{ij} \left| W_{i,j}^{c,h} - \frac{W_{i-1,j+1}^{c,h} + W_{i+1,j-1}^{c,h}}{2} \right| - \left| W_{i,j}^{c,h} - \frac{W_{i-1,j-1}^{c,h} + W_{i+1,j+1}^{c,h}}{2} \right|}{|W^{c,h}| \sum_{ij}},$$

# Extended Abstract Track

where  $\overline{|W|}$  denotes the average absolute value of the whole matrix and the sum is taken over all elements not on the edge of the matrix.  $\sum_{ij}$  in the denominator denotes the number of such elements. The translational symmetry shown in the plots is the symmetry summed over all circuits and heads (where the number of heads is denoted by  $N_h$ ):

$$\mathcal{S}_T = \sum_{c \in \text{OV, QK}} \sum_{h=1}^{N_h} \mathcal{S}_T^{c,h}$$

If the circuit is perfectly translationally symmetric, the irregularity of lines parallel to the diagonal will be 0. List sorting is translationally symmetric away from the vocabulary boundary, as a sorting algorithm only depends on the difference between elements within the list, not on their magnitude. We expect this to manifest in the circuits by having lines parallel to the diagonal be fairly uniform.

**Head Overlap** is introduced to measure the overlap between circuits of two heads  $h$  and  $h'$ . We take the sum of the absolute difference between elements in the heads normalized by the sum of the absolute value of the elements of the two matrices:

$$\mathcal{O}_{h,h'}^c = 1 - \frac{\sum_{ij} |W_{ij}^{c,h} - W_{ij}^{c,h'}|}{\sum_{ij} |W_{ij}^{c,h}| + |W_{ij}^{c,h'}|}.$$

The overlap shown in the plots are the mean overlaps of all OV-OV and QK-QK circuit combinations:

$$\mathcal{O} = \frac{1}{2} \frac{1}{N_h} \left( \sum_{c \in \text{OV, QK}} \sum_{h'=1}^{N_h} \sum_{h \neq h'}^{N_h} \mathcal{O}_{h,h'}^c \right)$$

## Appendix F. Varying the Model Architecture and Training

In this subsection, we study the impact of varying the model architecture and training such as the number of attention heads, and the use of LN and WD.

### F.1. 1-Head Model

Fig. 3 shows a single head transformer trained on our list-sorting task. As has been previously noted by [Bagiński and Kolly \(2023\)](#), a single head suffices for list sorting, and the model reaches 100% accuracy at step 133 (left panel of Fig. 3). After the model reaches 100% accuracy, the LLC keeps rising until step 672, with the circuit heat maps looking slightly less noisy and the translational symmetry rising strongly. Thereafter, the LLC drops together with The Circuit Rank as the off-diagonal stripes start to form at step 1985. These stripes become more pronounced towards the end of training (right panel of Fig. 3), but that doesn't seem to have an impact on the LLC. The OV and QK circuits seem to qualitatively function the same in all the panels. This is not surprising, as this model can't undergo head specialization.

The LLC has been calculated with inverse temperature  $n\beta = 512 / \ln 512 \approx 82$ , step size  $\epsilon = 10^{-4}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 2000$ .

## Extended Abstract Track

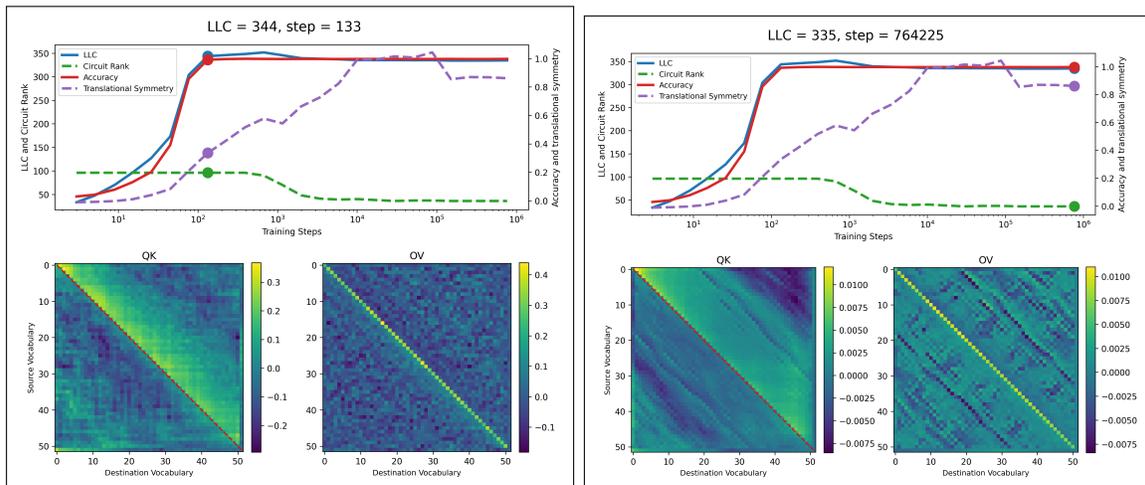


Figure 3: **1-head** list sorting model. As the LLC and Circuit Rank drop, off-diagonal patterns appear in the OV circuit (right) that are not present as the model learns how to sort well (left).

## F.2. Baseline 2-Head Model

In Fig. 4 we show the intermediate steps for the evolution of the baseline 2-head model. As previously mentioned, the heads **learns to sort early in the training** at step 133 (upper left of Fig. 4). A few hundred training steps later, the **LLC peaks and vocabulary-splitting begins to emerge** in the OV circuit, although the head overlap is still quite large and the QK circuit doesn’t show any distinct partitioning yet.

Next, the **LLC begins to drop while the heads clearly specialize into vocabulary-splitting**, with increasingly non-overlapping regions. We show a snapshot at step 3410 (lower left in Fig. 4), as the LLC stabilizes. The OV circuits of both heads are copying distinct regions. Accordingly, the QK circuit has also developed differentiated patterns along these regions. We expect and observe the attention pattern along dominant destination vocabulary regions to smoothly decrease from left to right, above the diagonal.

These features grow more distinct at the end of training, characterized by the appearance of off-diagonal patterns in the OV circuits. These begin forming after the LLC decrease and finish crystallizing after around 87k training steps. This is well captured by the translational symmetry measure.

We use a single layer attention only transformer model trained to sort lists of numbers. The baseline model architecture includes a residual stream size of 96, 2 attention heads with head dimension of 48, and Layer Normalization (LN) (Ba et al., 2016). The model is trained with Weight Decay (WD) set to 0.005 using the Adam optimizer (Loshchilov and Hutter, 2019) with a learning rate of  $10^{-3}$ , a dataset size of 150000 with a batch size of 512 and a cross-entropy loss function. The architecture is implemented using TransformerLens v.2.1.0 (Nanda and Bloom, 2022). The machinery used for training the models is NVIDIA RTX-4090.

# Extended Abstract Track

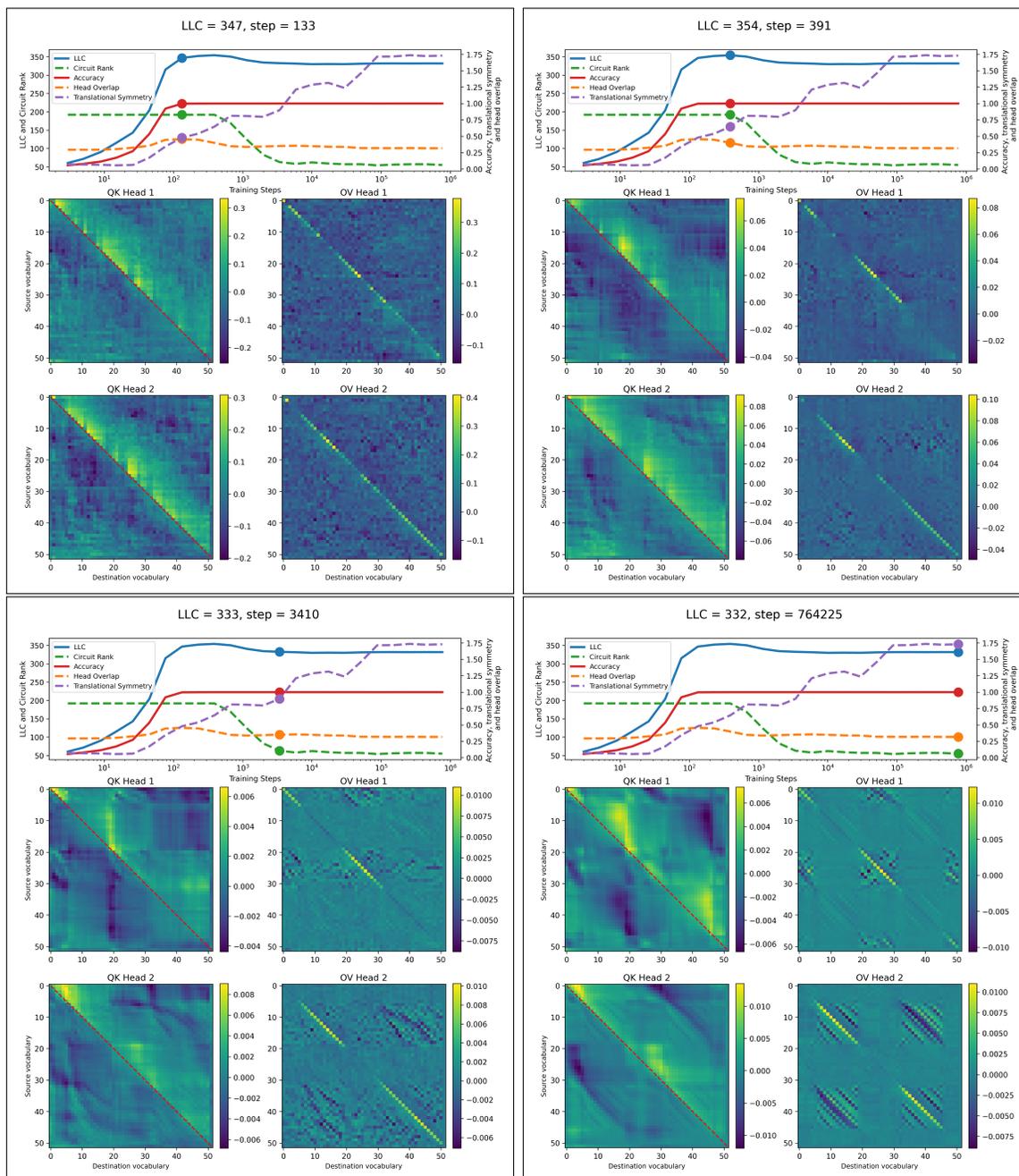


Figure 4: **Baseline 2-head model** as the model reaches 100% accuracy at step 133 (top left panel), peak LLC at step 391 (top right panel), after the LLC has dropped at step 3410 (bottom left), and at the end of training at step 764225 (lower right).

## Extended Abstract Track

**F.3. 3-Head Model**

As shown in the first row of Fig. 5, the 3-head model learns to sort at 100% accuracy at step 133, where all heads attend to and copy overlapping vocabulary regions. At peak LLC (2nd row of Fig. 5) we see first signs of vocabulary-splitting head specialization. As the LLC drops, the overlap between their vocabulary regions decreases, resulting in contiguous regions split across three heads, with head 3 covering only a small region (3rd row of Fig. 5). The QK circuits also display differentiated patterns, which upon closer inspection match the active vocabulary regions of the OV circuits. So far, the developmental stages of this model, match those of the baseline 2-head model.

As the evolution continues, around training step 5859 (not shown) the OV circuit of head 3 specializes into an "anti-state", seemingly suppressing the contributions from the other two heads, which behave like in the baseline 2-head model. We identify the state of head 3 to be copy-suppression, as discussed by McDougall et al. (2023). As the transition occurs, the QK circuit of head 3 also switches to uniform diagonal patterns, not differentiating any vocabulary regions anymore. This transition is not captured by any of our measures. This specialization can be seen at the end of training (4th row of Fig. 5).

The LLC has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 10^{-4}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 60000$ .

**F.4. 4-Head Model**

Similar to the other models, the 4-head model also learns to sort with 100% accuracy at step 133 (1st row of Fig. 6). As the LLC decreases, heads begin to specialize with concurrent vocabulary-splitting and copy-suppression appearing in heads 1,3,4 and head 2 respectively (2nd row of Fig. 6). The vocabulary regions are split unevenly, with head 4 covering only a very small region of the vocabulary.

This changes later in the training, after around 87k training steps (3rd row of Fig. 6), with heads 3 and 4 now copying similar vocabulary regions and displaying differentiated attention patterns in the QK circuits. Another transition is visible after 150k training steps (4th row of Fig. 6), where head 3 grows to attend and copy the entire vocabulary range. It seems to be suppressing the copy-suppression in head 3. This last transition is captured by a small drop in the LLC.

The model remains largely unchanged after this point, until the end of training (5th row of Fig. 6), as is seen from the measures remaining fairly constant.

The LLC has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 10^{-4}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 2000$ .

**F.5. Baseline 2-head Model without LN**

Removing LN from the baseline 2-head model causes a dramatic change to the training dynamics, as shown in Fig. 7. Early in training, at steps 71-348 (top row) the model goes through a transition in which The Circuit Rank drops dramatically, the LLC has a sharp increase and the head overlap drops. During this transition, the circuits of the model form a very regular dipole-like pattern.

This dipole-like pattern starts breaking at steps 18298-27194 (middle row) as the LLC peaks, with a formation of the stripe-like patterns parallel to the diagonal in the OV circuit.

# Extended Abstract Track

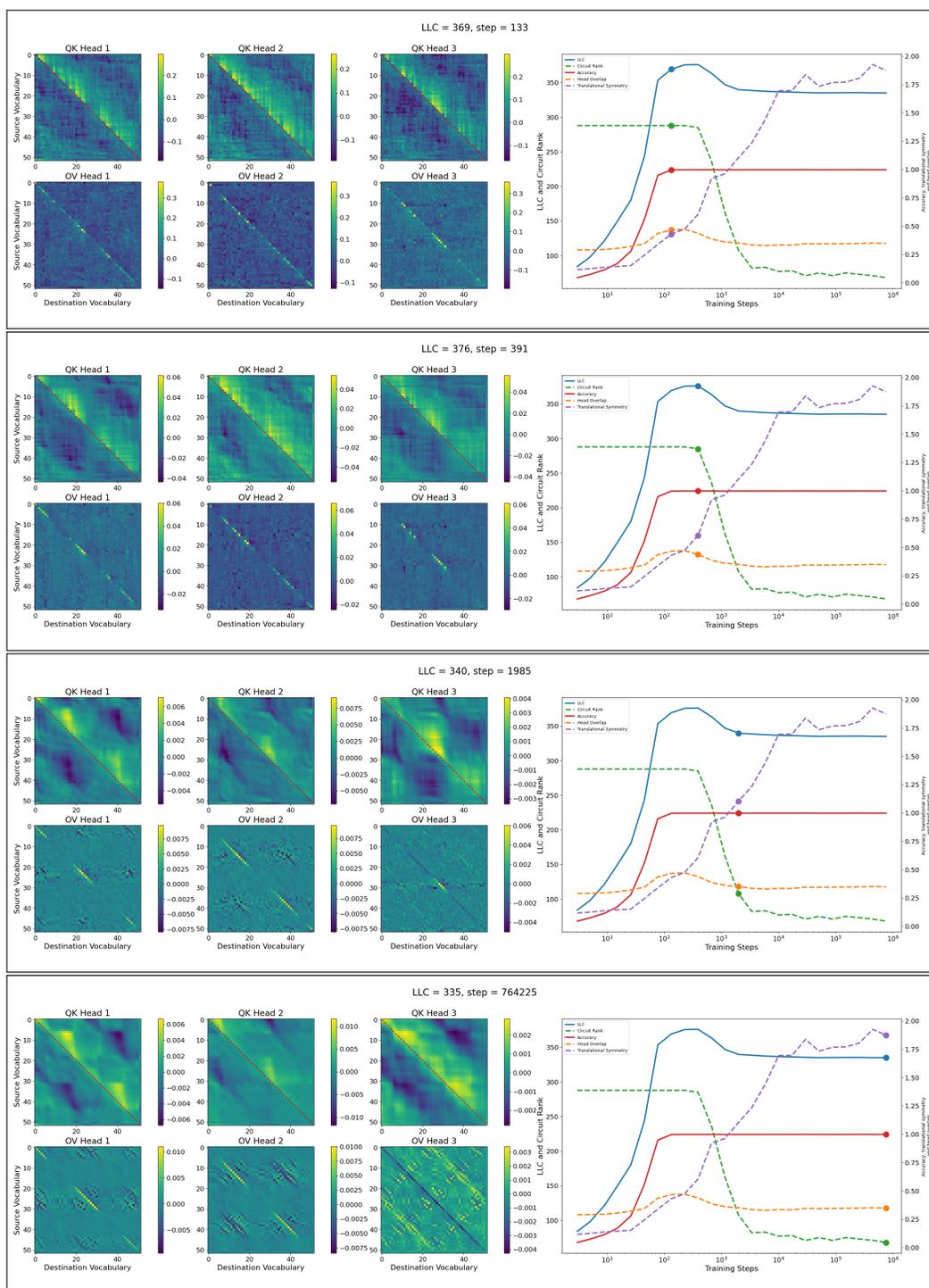


Figure 5: **3-head model** as the model learns how to sort (1st row), at LLC peak (2nd row), three-way vocabulary-splitting after LLC decrease (3rd row) and head 3 performing copy-suppression (4th row).

# Extended Abstract Track

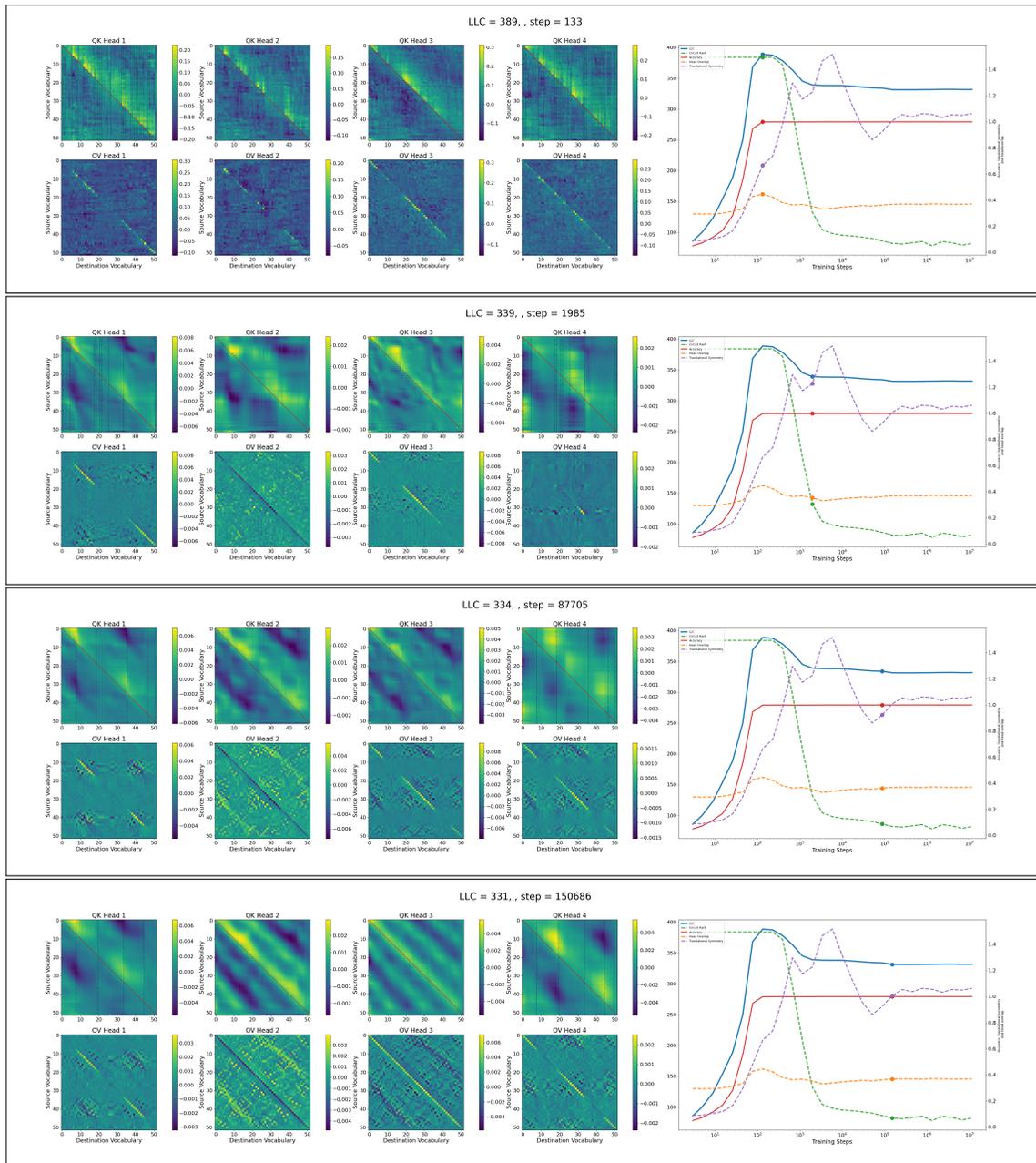


Figure 6: **4-head model** as the model learns how to sort (1st row), as the LLC decreases and heads specialize differently (2nd row), as heads 3 and 4 cover the same vocabulary regions (3rd row), as head 3 covers the entire range (4th row), and at the end of training (5th row).

# Extended Abstract Track

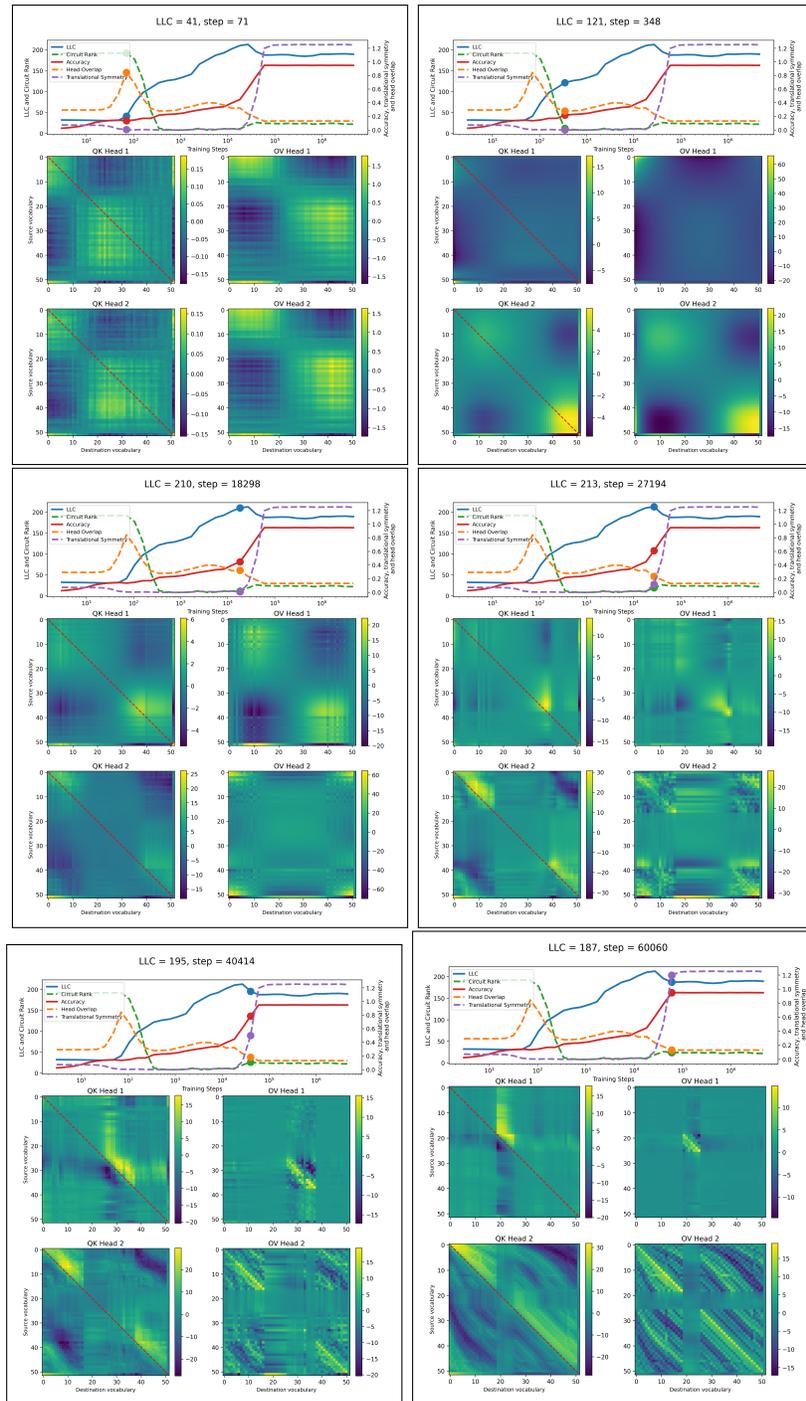


Figure 7: **Baseline 2-head model trained without LN** as the model simplifies but performs poorly (1st row), as relevant structure develops and performance improves rapidly (2nd row), as vocabulary-splitting appears before and after LLC decrease (3rd row).

# Extended Abstract Track

The QK circuits cover the regions determined by the OV circuit, similar to what we have seen in the other models. This structure formation stabilizes as the LLC drops (bottom row), which is also tracked by a dramatic increase in the translational symmetry. The model never reaches 100% accuracy on list sorting, and the accuracy does not flat-line until step 60060, after which all the measures are stable. The LLC seems to capture the development of this model very well.

The LLC has been calculated with inverse temperature  $n\beta = 26$ , step size  $\epsilon = 10^{-6}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 3$  and  $n_{\text{draws}} = n_{\text{burnin}} = 100000$ .

## F.6. Baseline 2-head Model without WD

As seen in Fig. 8, the model without WD still learns to sort at 100% at step 133. The OV and QK circuits seem more noisy, and there is no drop in the Circuit Rank. The LLC still has a large drop between steps 1985 and 10066 during which the heads specialize into splitting the vocabulary. This specialization is clearer for tokens smaller than 20 in the QK and OV circuits, less so for larger vocabulary tokens.

The LLC has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 3 \times 10^{-6}$ , localization term  $\gamma = 56$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 65000$ .

## F.7. Baseline 2-head Model without LN and WD

Fig. 9 shows the evolution of our measures and the circuits for the baseline 2-head model without both LN and WD. Compared to the baseline model, it learns to sort at 100% accuracy somewhat later, at step 391, but considerably faster than the baseline model without LN. The model seems to go via dipole-like circuits around step 45 as the head overlap peaks, very similar to step 71 of the baseline model without LN (compare the top left panels of Figs. 7 and 9). Instead of going via the low Circuit Rank dipole phase, however, the model instead develops circuits that are capable of sorting, while still retaining some of the dipole-like patterns at step 391. This happens at the same time as the LLC peaks.

After this, the LLC drops, and the dipole like pattern gives way to patterns resembling the baseline 2-head model, with partial vocabulary-splitting head specialization in both QK and OV for vocabulary below around 20. We speculate that the reason why the presence of WD causes a worse performance is that it pushes the circuits into simpler low-rank dipole-like patterns instead of learning to sort.

The LLC has been calculated with inverse temperature  $n\beta = 30$ , step size  $\epsilon = 10^{-6}$ , localization term  $\gamma = 56$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 40000$ .

## Appendix G. Varying the dataset

In this subsection, we study the impact of varying aspects of the training data, such as the size of the vocabulary, the length of the list and the presence of perturbations in the data set.

### G.1. Baseline 2-head Model with Vocabulary Size Increased to 202

Increasing the vocabulary size to 202 produces the training dynamics shown in Fig. 10. The model reaches 97% accuracy around step 635, which coincides with the LLC peak.

# Extended Abstract Track

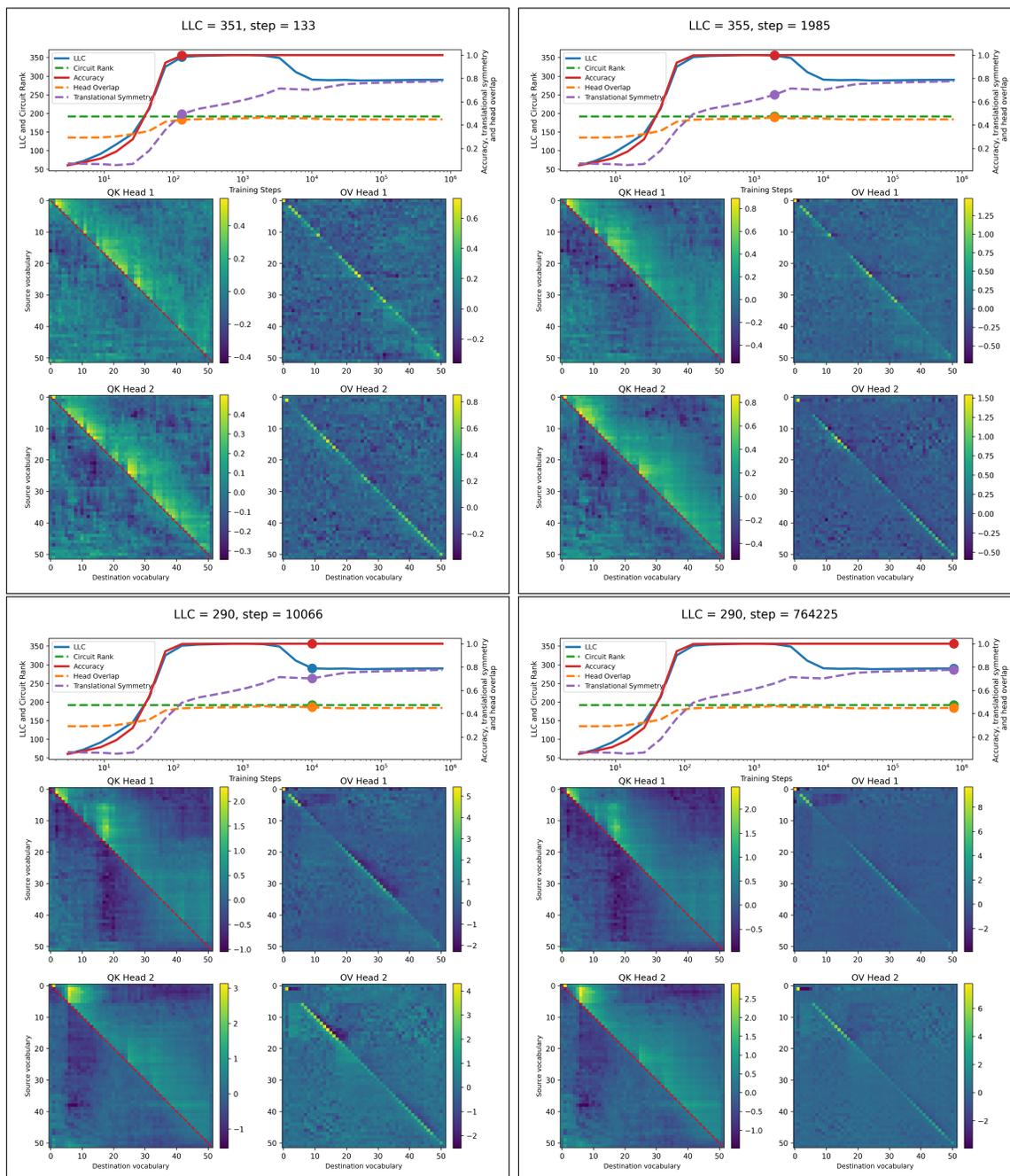


Figure 8: **Baseline 2-head model trained without WD**, as the model learns how to sort (upper left), as the LLC is at its peak (upper right), after the LLC drop (lower left) and at the end of training (lower right).

# Extended Abstract Track

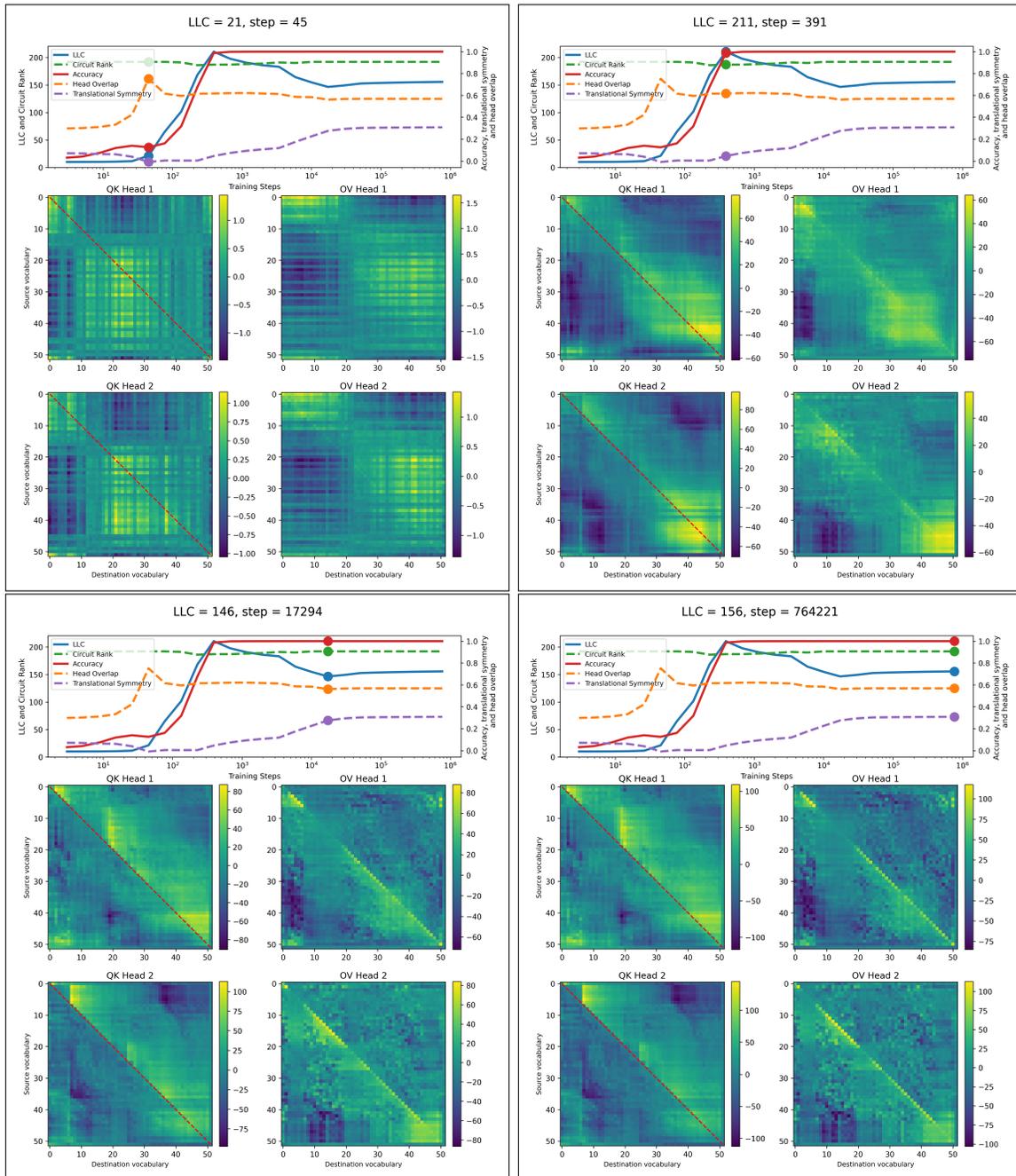


Figure 9: **Baseline 2-head model without LN and WD** as head overlap peaks (upper left), as accuracy is high and LLC peaks (upper right), after LLC drop (lower left) and at the end of training (lower right).

# Extended Abstract Track

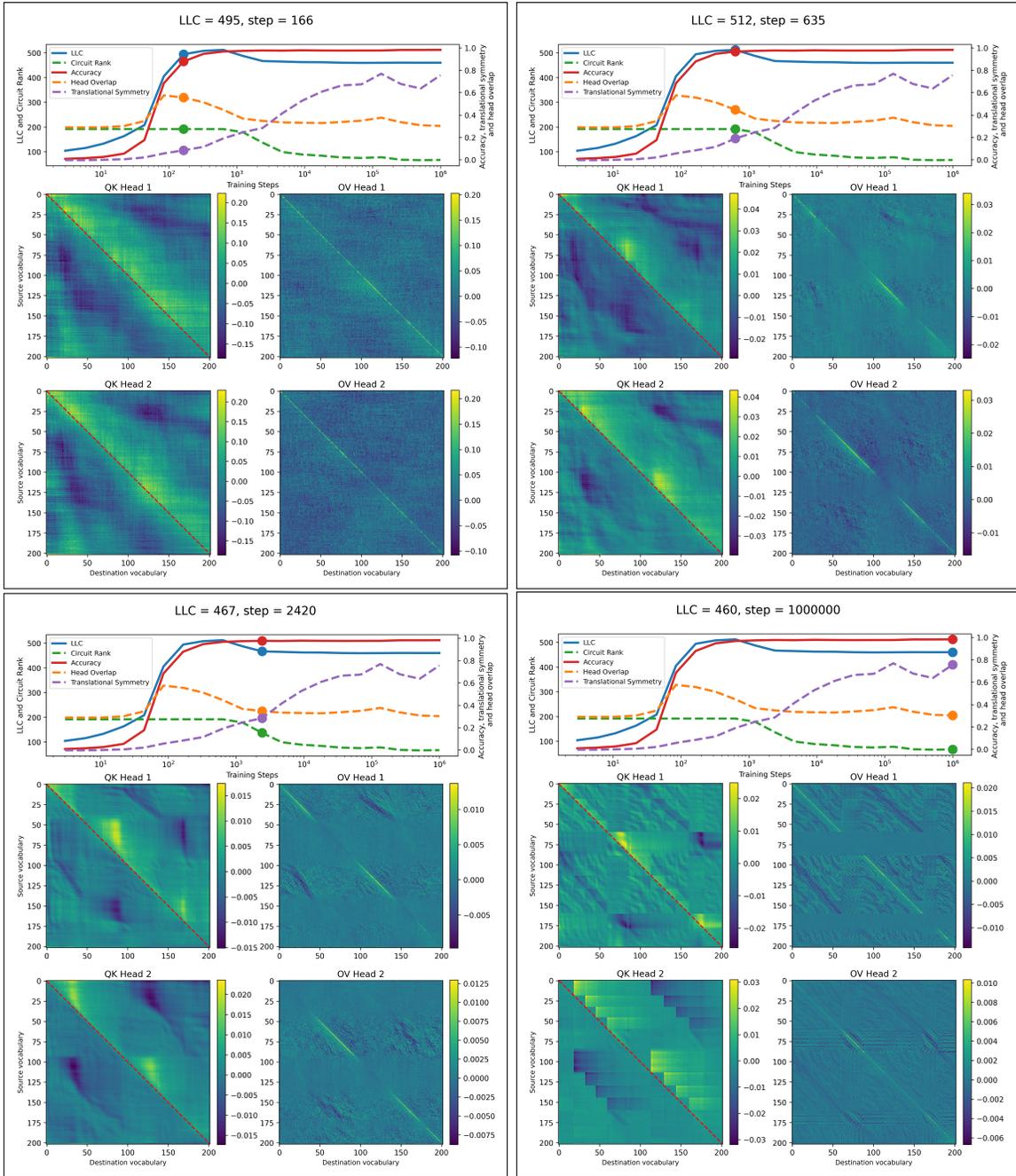


Figure 10: Baseline 2-head model with **vocabulary size increased to 202**, we find similar developmental stages as in the baseline model. **Vocabulary region size increases.**

# Extended Abstract Track

Thereafter, the LLC and the Circuit Rank drop as the heads specialize into vocabulary-splitting. At the end of training, the model develops a square-like pattern in the QK circuit, which doesn't always correspond to a vocabulary region boundary. This last transition is accompanied by a small drop in The Circuit Rank, but no drop in the LLC.

The LLC has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 10^{-3}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 2000$ .

## G.2. Baseline 2-head Model with List Length Increased to 20

Increasing the list length to 20 yields the training dynamics shown in Fig. 11. Here, the LLC peaks as the model reaches 100% accuracy at step 166, and then drops as the heads specialize into contiguous regions of parameter space and The Circuit Rank drops. We note that compared to the baseline 2-head model, the larger list length leads to a larger number of regions, distributed in a periodic pattern.

The LLC has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 3 \times 10^{-6}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 70000$ .

## G.3. Baseline 2-head Model with Perturbed Dataset

We perturb the data by iterating through the dataset once, and swapping neighboring elements in the sorted list with probability  $40\%/(n_{i+1} - n_i)$ , where  $n_i$  is the value of the list element  $i$ . Since the probability of neighboring elements swapping is always less than 50%, we believe that the optimal strategy still should be to sort the list ignoring the perturbations. The perturbations do, however, have a severe impact on the training dynamics, as shown in Fig. 12.

We don't observe any drop in the LLC, even though The Circuit Rank does drop. The heads don't specialize into vocabulary-splitting modes, but the OV circuits rather settle into what looks like opposites of each other. It looks like head 1 does copy suppression and head 2 does copying, whereas the QK circuits behave very differently from what we have seen in the other models.

The accuracy has been computed on non-perturbed data, and increases throughout training, reaching 98% at the end of training.

The LLC has been calculated with inverse temperature  $n\beta = 512/\ln 512 \approx 82$ , step size  $\epsilon = 10^{-6}$ , localization term  $\gamma = 32$ ,  $n_{\text{chains}} = 4$  and  $n_{\text{draws}} = n_{\text{burnin}} = 200000$ .

## Appendix H. Compute

For the experiments, we rented RTX-4090s on vast.ai. We spent a total of \$200, giving about 650 GPU hours. Some experiments had to be re-run and some were not used, and we estimate that about 70% of the compute went into the results included in the paper. The cost of the LLC estimation is proportional to

$$n = n_{\text{chains}} \times (n_{\text{draws}} + n_{\text{burnin}}) = 2n_{\text{chains}} \times n_{\text{draws}}$$

For each model, the LLC has been computed at 24 snapshots, and for each model parameter sweeps have a total cost of approximately  $n \approx 2.6\text{M}$ . Adding LLC cost of all the experiments used in the paper together we get  $n \approx 131\text{M}$ . Of the compute used in the paper, we

# Extended Abstract Track

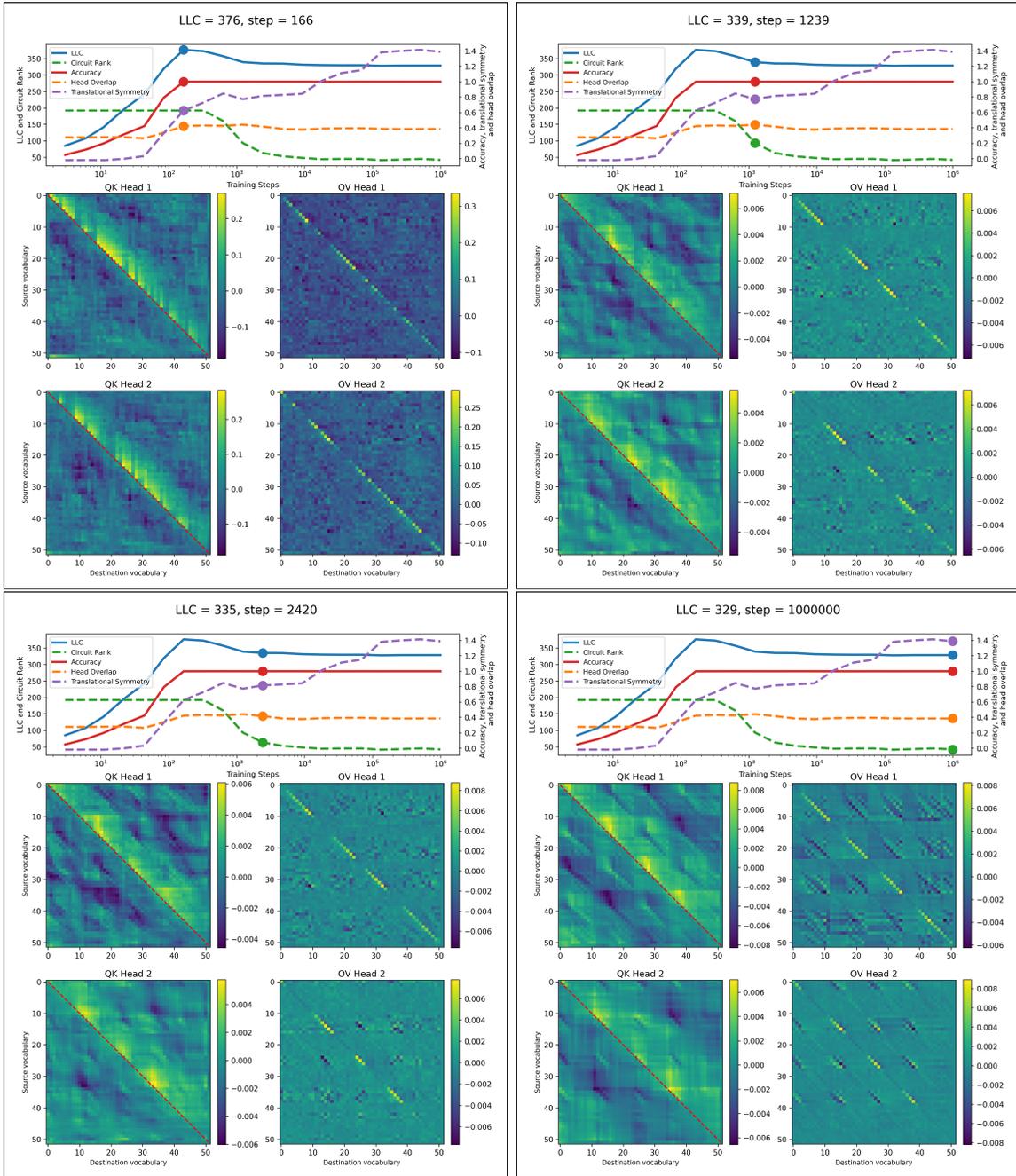


Figure 11: Baseline 2-head model with **list length increased to 20**, we find similar developmental stages as in the baseline model. **Number of vocabulary regions increases.**

# Extended Abstract Track

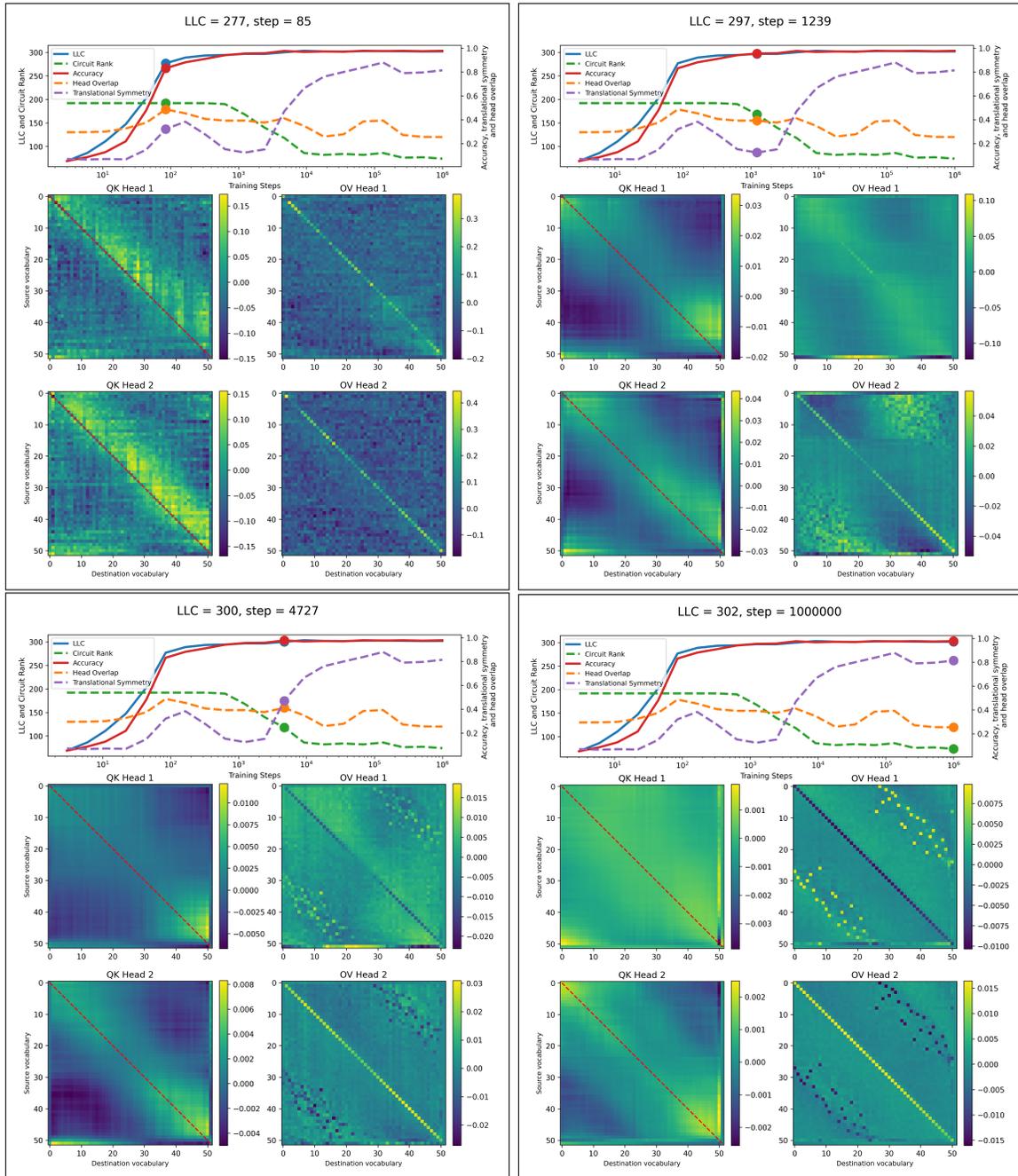


Figure 12: Baseline 2-head model with **perturbed training dataset** shows different developmental stages and it is the only 2-head model where we observe **copy-suppression**.

# Extended Abstract Track

estimate that we spent around 30% on training the models, and 70% on LLC estimation. The computation cost of an experiment, including model training, LLC hyperparameter scan and LLC estimation can be estimated as

$$\left(0.03 + 0.7 \times \frac{2.6\text{M} + 24 \times n_{\text{chains}} \times (n_{\text{draws}} + n_{\text{burnin}})}{131\text{M}}\right) \times 0.7 \times 650 \text{ GPU hours},$$

where  $n_{\text{chains}}$ ,  $n_{\text{draws}}$  and  $n_{\text{burnin}}$  is stated for every experiment.