TRANSFORMER MECHANISMS MIMIC FRONTOSTRI ATAL GATING OPERATIONS WHEN TRAINED ON HU MAN WORKING MEMORY TASKS

Anonymous authors

Paper under double-blind review

ABSTRACT

The Transformer neural network architecture has seen success on a wide variety of tasks that appear to require *executive function* – the ability to represent, coordinate, and manage multiple subtasks. In cognitive neuroscience, executive function is thought to rely on sophisticated frontostriatal mechanisms for selective gating, which enable role-addressable updating- and later readout- of information to and from distinct "addresses" of memory, in the form of clusters of neurons. However, Transformer models have no such mechanisms intentionally built-in. It is thus an open question how Transformers solve such tasks, and whether the mechanisms that emerge to help them to do so resemble the gating mechanisms in the human brain. In this work, we analyze the mechanisms that emerge within a vanilla attention-only Transformer when trained on a task from computational cognitive neuroscience explicitly designed to place demands on working memory gating. We find that the self-attention mechanism within the Transformer develops input and output gating mechanisms, particularly when task demands require them. These gating mechanisms mirror those incorporated into earlier biologically-inspired architectures and mimic those in human studies. When learned effectively, these gating strategies support enhanced generalization and increase the models' effective capacity to store and access multiple items in memory. Despite not having memory limits, we also find that storing and accessing multiple items requires an efficient gating policy, resembling the constraints found in frontostriatal models. These results suggest opportunities for future research on computational similarities between modern AI architectures and models of the human brain.

034

006

007

008 009 010

011 012 013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

033

036

038

1 INTRODUCTION

The Transformer architecture Vaswani et al. (2017) has recently become the dominant neural network model in artificial intelligence. Unlike some earlier AI architectures, which were inspired (albeit loosely) from human processing of language (Hochreiter, 1997) or vision (LeCun et al., 2015), Transformers have no mechanisms designed overtly to resemble the human brain. It thus remains an open question what, if any, similarities exist, and whether there are opportunities for theory or insights from AI and neuroscience to mutually inform one another (McGrath et al., 2024).

In this work, we focus specifically on working memory management via *gating mechanisms*. Gating mechanisms are responsible for multiple distinct aspects of working memory management, includ-ing determining which items to store and retrieve, when, and from where. In humans, there is strong evidence that such mechanisms are essential for tasks that require *executive function*, i.e., the ability to manage competing demands from multiple tasks, stimuli, and responses and to coordinate their execution (Frank & Badre, 2012; Badre & Frank, 2012; Chatham et al., 2014; Rac-Lubashevsky & Kessler, 2016; Rac-Lubashevsky & Frank, 2021). Although some Transformer variants have additional built-in structure for memory (Dai et al., 2019; Burtsev et al., 2020; Wang et al., 2019), the architectures which currently dominate modern AI systems are "vanilla" (Brown et al., 2020; Touvron et al., 2023), lacking specialized components to support this type of control.

054 We thus investigate whether, and under what conditions, such structure can emerge as a result of 055 training. We train vanilla Transformer models on a task from cognitive neuroscience that was de-056 signed designed to investigate selective gating and working memory in humans (O'Reilly & Frank, 057 2006; Rac-Lubashevsky & Frank, 2021). We use recent techniques from mechanistic interpretability 058 (Olah, 2022; Nanda & Bloom, 2022) to expose the mechanism that the Transformer uses to perform the task. We find that, as a result of training, the self-attention mechanism specializes in a way that resembles existing models of input-output gating (§4.1), but that these mechanisms only arise when 060 the training task places specific demands on gating that mimics biological networks (§4.2). We fur-061 ther find that when such mechanisms do arise, they are predictive of better task performance and 062 of generalization to changes in the input distribution and task demands (§4.3), improving effective 063 working memory capacity (§4.4). Our findings highlight the importance of considering the emer-064 gent mechanisms that result from training in addition to the innate architectural mechanisms when 065 drawing comparisons between AI systems and human cognitive processes. 066

067

2 BACKGROUND AND HYPOTHESES

068 069

There is strong evidence that working memory in human brains makes use of a gating mechanism to read, write, and maintain information required to carry out complex tasks (Rac-Lubashevsky & 071 Kessler, 2016; Rac-Lubashevsky & Frank, 2021; Bhandari & Badre, 2018; Badre & Frank, 2012; 072 Chatham et al., 2014). Sophisticated gating mechanisms of the type implemented in biological 073 neural networks contain at least three important components. First, *input gating* controls whether 074 or not given information is stored in memory, and if stored, determines the "address" (population of 075 neurons) to which it should be written. Second, output gating determines when and what information 076 to read out of memory to inform a subsequent decision, such as to produce a response to a task. 077 Finally, working memory is role addressable, meaning that items can be bound to a learned taskdependent context (i.e. role) when stored and accessed in working memory. For example, in listening 079 to a story for the first time, people will not just remember individual entities ("Andrea", "Chicago") but rather can associate them with specific roles such as lives_in ("Andrea", "Chicago").

081 In humans and other animals, these operations are supported by corticostriatal circuits in which 082 isolated clusters of prefrontal neurons (or *stripes*) are used to represent distinct addresses in memory 083 that can be updated or read out from via selective gating actions triggered by basal ganglia and 084 thalamus (O'Reilly & Frank, 2006; Frank & Badre, 2012; Kriete et al., 2013; Calderon et al., 2022; 085 Soni & Frank, 2024). These stripes can also serve as latent roles that condition how to interpret content within them. When learning effective gating policies, these models afford functions such as variable binding and indirection that support rapid generalization to new situations (O'Reilly & 087 Frank, 2006; Frank & Badre, 2012; Collins & Frank, 2013; Kriete et al., 2013; Bhandari & Badre, 880 2018). 089

In principle, Transformers are good candidates for learning such gating behavior. Transformers' na tive self-attention mechanism consists of attention heads which are arguably functionally analogous
 to frontostriatal stripes. The decomposition of these heads into distinct keys, queries, and values
 (see Appendix A.1) means that the Transformer can in principle learn to differentiate reading and
 writing operations in a role- and context-dependent way across its multiple heads. However, whether
 Transformers will use their self-attention to implement such a mechanism is an open question, especially in cases when it is possible to fit the training data using more heuristic and less generalizable
 solutions.

We thus consider two hypotheses. First, we investigate whether, and under which data distributions, the Transformers use their attention heads to learn effective gating strategies. We find that the key vectors form addresses analogous to the PFC "stripes" (neural populations that support variable binding in memory). The learned key construction determines the address to store an item and is thus analogous to input gating. Conversely, the query vectors determine which addresses are accessed, and are thus analogous to output gating.

Second, we investigate whether adopting such strategies will facilitate rapid learning and generalization in working memory tasks the way it has been show to in humans. Specifically, human studies
have shown that working memory capacity is not limited by the number of items one can maintain
but rather by their effective gating strategies in frontostriatal circuits (Vogel et al., 2005; McNab &
Klingberg, 2008; Baier et al., 2010). Theoretical work has shown that capacity limits in these circuits

do not stem from a limitation in the number of available neural populations, but rather result from a credit assignment problem that arises when learning to manage multiple items in memory (Soni & Frank, 2024; Todd et al., 2009). These limitations can thus be partially mitigated by learning to reuse effective gating policies (Soni & Frank, 2024). Because these limitations are computational rather than anatomical (i.e., not driven by the number of neurons/populations available), we hypothesize that they would also manifest in Transformers, even though they have no inherent memory demands at all (since all information is available in the context window).

115 116

117

3 EXPERIMENTAL DESIGN

118 3.1 TASKS

119 Reference-Back 2 Task: We use a variant of a task from cognitive neuroscience known as the 120 "reference-back 2" task (Rac-Lubashevsky & Frank, 2021). This is one among a number of task 121 designs inspired by frontostriatal modeling work (O'Reilly & Frank, 2006; Soni & Frank, 2024) 122 which requires selective updating and accessing of information in a role-addressable manner. In the 123 reference-back paradigm, symbols are viewed one at a time with associated roles, and the participant 124 must determine whether the current symbol is the same or different as that stored in memory for a 125 given role. For example, a sequence might contain letters (role 1) and numbers (role 2), each of 126 which occurs along side an update instruction which is either Store or Ignore. For each symbol 127 in the sequence, the participant must do two things: 1) make a same/different judgment based on 128 whether the current symbol matches the previously-stored symbol for that role, and 2) if the update instruction is Store, update the symbol associated with the associated role. See Figure 1 for an 129 example. We create a modified text-based version of the reference-back 2 task. In our design, roles 130 are denoted explicitly using special tokens (i.e., either Reg0 or Reg1 for the two-role version) 131 indicating the role (or "register") to which the symbol should be bound. See Appendix A.2 for more 132 details about our task implementation. 133



Figure 1: **Task** The reference-back-2 task requires making same vs. different judgments for each symbol in a sequence by comparing against a previously-shown symbol. See text for description of the task. In the above example, there are two roles, blue and red. The register state (shown along the top and connected by dotted lines) is assumed to be latent in the model; i.e., it is not provided as input.

149 150

157

134 135

136 137

143 144

Split-Set Control: The computational advantages of role-addressable gating in frontostriatal networks (relative to other recurrent neural networks) are particularly evident when any symbol can be assigned to any register, so that the networks have to learn to assign them separable addresses (O'Reilly & Frank, 2006). To test the hypothesis that gating mechanisms only emerge in response to such task demands, we create a control condition in which the registers are associated with disjoint sets of symbols. In this task, it is possible to succeed without role-addressable gating, since symbols never need to be decoupled from their associated registers.

Ignore-Integrated and Ignore-Separated Controls (a.k.a Split-Set Control): In the same vein,
we include a further simplified task variant which is designed to require selective input gating but not
output gating. This condition, along with the split-set control, have been previously used to show
that the advantage of frontostriatal gating networks relative to other recurrent networks is largely
reduced in such scenarios (O'Reilly & Frank, 2006).

Specifically, in the *ignore separated* condition (referred to as split-set control earlier), each register sees a disjoint set of symbols. Additionally, symbols falling under the Ignore update instruction are disjoint from those that are stored, meaning the model does not need to learn to differentiate Store and Ignore in a meaningful way. In contrast, in the *ignore integrated* condition, symbols under an Ignore instruction will be in distribution with the chosen register. In other words, the model must distinguish between a Store and Ignore instruction and thus learn input gating, but does not have to learn an output gating policy because both registers have mutually exclusive symbols.

170 171

3.2 MODELS

We train small, attention-only, decoder-only Transformer models from scratch on our task. Our
models contain two decoder-only layers, each with two heads, and no multilayer perceptrons or
layer normalization, followed by a linear "unembed" layer. The models are trained on 100k training
data points for 60 epochs. See Appendix A.3 for additional details.

177 178 3.3 METRICS

Performance on Refback2 Task: We evaluated how well the model performs on the task on which it was trained using standard accuracy on the same vs. different prediction for each symbol in the sequence.

182

183 **Input and Output Gating:** The refback2 task is assumed to benefit from gating mechanisms, 184 but success on the task is not in and of itself diagnostic of having learned the gating mechanisms. 185 To develop an intrinsic measure of the input and output gating mechanisms, we use path-patching (Wang et al., 2022; Goldowsky-Dill et al., 2023), a generalization of causal mediation analysis (Pearl, 2001) that allows us to determine which components of a neural model (e.g., attention heads) 187 work together in order to produce observed behavior on a task. Path patching involves designing 188 a minimal pair of inputs, a "clean" input and a "corrupted" and then finding specific components 189 of a model which fully account for the difference in the model's output between the two cases. By 190 "patching in" only these components, the model can be made to behave as though it is seeing the 191 corrupted input even when it is in fact seeing the clean input. See Appendix A.4 for more details. 192

We use path patching as a measure of input gating in the following way. For input gating, we design a minimal pair of inputs in which the corrupted copy includes Ignore in a place where the clean copy contained a Store, or vice-versa. We then use path patching to identify an incisive edit that can be made to the weights of the model in order to prevent the model from storing ("gating in") a given symbol. We report the accuracy of the input gating mechanism as the percentage of inputs on which this edit to the weights changes the models behavior in the expected way.

Analogously for output gating, we design minimal pairs which we expect to yield differences in an
 output gating mechanism, assuming one exists and is functioning correctly. Specifically, our clean
 and corrupted sequences differ in the register associated with one of the symbols, which should trigger a difference in what information is read ("gated out") in order to make a final same vs. difference
 judgment. Again, we report the accuracy of the output gating mechanism as the percentage of cases
 on which it is possible to make such an edit and produce the desired effect.

205 Our description of results in Section 4 elaborates on both this and the input gating metrics and their 206 interpretation.

207 208

4 Results

- 209
- 210 211

4.1 KEY AND QUERY VECTORS SPECIALIZE FOR INPUT AND OUTPUT GATING

We find that Transformers implement a role-addressable gating mechanism in which key vectors control input gating and query vectors control output gating (see Appendix A.1 for summary of key, query, value attention). Specifically, after training, the key for Sym_i (e.g., at tokens 2, 6, 10, and 14 in Fig. 2a) represents the combination of the update instruction, register, and symbol for position *i*. A query's ability to address this position depends on whether the represented tuple contains a Store

216	Clean sequence	Attention heatmap	Prediction (idx 15)
217	a) 0 1 2 3 3 4 5 5 6 5 7 8 5 9 10 11 11 12 13 14 5 5 15 14 5 15 14 15 15 15 15 15 15 15 15 15 15 15 15 15	Optione Reg0 Sym5 3 diff 4 store 5 Reg1 6 ym3 7 diff 8 store 9 Reg1 10 diff 11 lignore 12 lignore 13 lignore 14 diff 10 diff <th10 diff<="" th=""> 10 diff 10</th10>	same
219	Corrupted sequence (with minimal pair patched to target indices)	Attention heatmap after patch	
220	$\textbf{b} \begin{bmatrix} 0 & 1 & 2\\ \text{store} & \text{Reg0} & \text{Sym5} & \text{diff} & \text{store} & \text{Reg1} & \text{Sym3} & \text{diff} & \text{spore} & \text{Reg1} & \text{spore} & \text$	Open Rego Sym3 3diff store Store Sym3 7diff store Reg1 100 111 122 133 144 100 100 111 112 133 144 100	different
221 222		$ \underbrace{ \left[\begin{array}{c} 0 \\ \text{store} \end{array} \right] \left[\begin{array}{c} 2 \\ \text{gym5} \end{array} \right] \left[\begin{array}{c} 3 \\ \text{diff} \end{array} \right] \left[\begin{array}{c} 4 \\ \text{store} \end{array} \right] \left[\begin{array}{c} 6 \\ \text{Reg1} \end{array} \right] \left[\begin{array}{c} 7 \\ \text{store} \end{array} \right] \left[\begin{array}{c} 8 \\ \text{Reg1} \end{array} \right] \left[\begin{array}{c} 1 \\ \text{Store} \end{array} \right] \left[\begin{array}{c} 1 \end{array} \left[\begin{array}$	different
223	$\textbf{d)} \bigcup_{\text{store}} \frac{1}{\text{Rego}} \frac{2}{\text{SymS}} \frac{3}{\text{diff}} \frac{4}{\text{store}} \frac{5}{\text{Rego}} \frac{6}{\text{Sym3}} \frac{7}{\text{diff}} \frac{8}{\text{store}} \frac{9}{\text{Reg1}} \frac{10}{\text{sym4}} \frac{11}{\text{diff}} \frac{12}{\text{gnore}} \frac{13}{\text{Rego}} \frac{14}{\text{Sym4}} 14$	0 store 1 Reg0 Sym5 3 diff 4 store 8 Reg1 Sym3 7 diff 8 store 9 Reg1 Sym4 1 diff 1 12 reg1 Sym4 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1	same
224	$\textbf{e} \underbrace{\begin{smallmatrix} 0 & 1 \\ store & Reg0 & Sym5 \\ \hline \textbf{Sym5} & \text{diff} & \frac{4}{store} & \frac{5}{Reg1} & \frac{6}{Sym3} & \frac{7}{diff} & \frac{8}{store} & \frac{9}{Reg1} & \frac{10}{Sym4} & \frac{11}{diff} & \frac{12}{gnore} & \frac{13}{sym4} & \frac{14}{sym4} & \frac{14}{sym4$	0 store 1 Reg0 Sym5 3 airff 4 store 7 Reg1 Sym3 7 diff 8 store 9 Reg1 Sym4 1 11 12 Sym4 1 12 Sym4 1 11 12 Sym4 1 1	different
226	$\textbf{f)} \underbrace{\begin{smallmatrix} 0 & 1 & 2 & 3 \\ \text{store} & \text{Reg0} & \text{Sym4} \\ \end{bmatrix}}_{\text{Sym4}} \underbrace{\begin{smallmatrix} 3 & 1 & 4 \\ \text{diff} & \text{store} & \text{Reg1} \\ \text{Sym3} & \text{diff} & \text{store} \\ \end{bmatrix}}_{\text{Sym4}} \underbrace{\begin{smallmatrix} 10 & 11 & 12 & 13 \\ \text{gnore} & \text{Reg0} \\ \text{Sym4} \\ Sym4$	Store Reg0 Syms 3 attr 4 store Reg1 Sym3 7 attr 8 store Reg1 50m4 1 attr 12 attr 14 attr 10 attr 14 attr 10 attr 14 attr 10 attr 10 attr 10 attr 11 attr 12 attr 14 attr 10 attr 14 attr 10 attr 14 attr 10 attr 14 attr 10 attr 10 attr 10 attr 10 attr 10 attr 11 attr 12 attr 10 attr	same

229

230

231

232

233

> Figure 2: Path Patching Examples. Model behavior across different path patching conditions. Attention is visualized as a shade of purple, with deeper shade corresponding to higher attention to that token. We create "corrupted" minimal pairs in which changing a token in the input (light blue) either changes the correct label at index 15 (examples b, c, e) or does not (d, f). We make small path-patching edits with the minimal pair to targeted network components (layer 1 keys for b, c, d, f; queries for e,f, see text). In all test examples, making the small patch successfully alters the model's prediction to align with the "corrupted" example, as expected.

- 234 235
- 236

or an Ignore. That is, key vectors representing an Ignore tuple receive very little attention (0.4% 237 of layer 1 attention averaged over test set), whereas those representing a Store tuple receive the 238 bulk of the attention (86.8%). We demonstrate the above narrative using path patching, described 239 in Section 3.3), and further below. Figure 2 shows a summary of the path patching experiments and 240 results. Our task is simple and does not contain noise, so in all cases, the intervention (i.e., patching 241 to the keys or queries) results in a 100% change in the model prediction in the expected direction. 242 Thus, for compactness, Figure 2 depicts the conditions but does not include quantitative results.

243 First, to investigate input gating, we create clean sequences sampled from our test set, and then 244 corrupt these sequences by switching a Store within tuple i to an Ignore. We path-patch only the 245 key vectors of *i*. We expect, if the key controls input gating, that patching these key vectors should 246 "block" attention to all of tuple *i*. An example attention pattern is in Fig. 2, examples a and b. We 247 find that the model's attention shifts away from the tuple accordingly in 100% of patched instances. 248 The presence of an Ignore or a Store within a tuple controls whether the key construction acts 249 as an open input gate or a closed input gate.

250 Key construction also depends on the *role* of the represented content; within our task, that means 251 whether Reg_i is Register 0 or Register 1. When making a same/different prediction, key vectors 252 representing a tuple that matches the target register receive most of the model's attention (92.5% of 253 total attention), while those that do not match are not attended to (3.3% of total attention). Again, we 254 use path-patching to determine that key construction encodes roles, this time perturbing the target 255 register rather than the update instruction (see Fig. 2, row c). The model's attention shifts away accordingly across every example in the test set. Note that the stored tuple must be modified; if the 256 same corruption is made earlier (as in row d), attention does not shift. This behavior shows that the 257 gating within self-attention is *role-addressable*; the registers within the task function as roles, and 258 are embedded within the key vectors as part of the representation. 259

260 Given that key vectors serve the role of addresses, query vectors in turn control which key vectors are accessed, through the final Q*K dot product in attention. Query construction thus performs the role 261 of output gating within Transformers. The query composition controls which addressable Symbol i 262 representations are attended to based on the identity of the target register. We again determine this 263 through a set of path-patching experiments in which we perturb the target register (Fig. 2, row e). 264

265 We find that patching to the query vector in such cases indeed causes the attention to shift from 266 the original stored tuple (74.1% of attention) to the stored tuple that matches the edited register, 267 resulting in a corresponding change in the final same/different judgment. Editing aspects of the target tuple other than target register has minimal effect on the query construction. No edits to the 268 query cause the model to attend to an Ignore tuple, further evidencing of output gating behavior-269 only content that has been made "addressable" can be accessed for a response. Furthermore, we 270 find that the target instruction and symbol do not factor into the query composition- changing them 271 through path-patching to the query does not affect attention. This is notable because the model 272 could employ other strategies for determining which tuples are eligible to be the stored tuple; e.g. 273 attending to all symbols to match if any of them are the same as the target symbol.

EMERGENCE OF GATING POLICIES DEPENDS ON TASK DEMANDS 4.2

Under which conditions to gating computations arise? One possibility is that they are a trivial consequence of the Transformer's architecture. After all, gating is sometimes expressed as a simple multiplicative operation, as in LSTMs (Hochreiter, 1997). And indeed, the key, query, and value vectors underlying the attention heads are combined via matrix multiplications. However, this alone does not ensure that the networks learn effective, role-addressable gating policies. We thus investigate whether the training task demands induce such gating policies by running several control experiments (§3.1) which simplify the task, removing the role addressability of gating requirements. We thus hypothesized that since this task is much easier to learn, the Transformer will learn an overly memorized, brittle solution and will not develop effective input and output gating strategies.



300 Figure 3: Split Set Control a) Ignore Symbols are not mutually exclusive with the registers. The models must learn to respond appropriately to store/ignore instructions (input gating). This control task separates the symbols shown to each register, reducing the need for the model to learn output 302 gating. b) Both registers and the Ignore instruction have mutually exclusive symbols. 303

As shown in Figure 3, networks trained on the split set control task do not perform well at either 305 the input or the output gating subtasks (see §3.3 for gating metrics). Moreover, models trained on 306 the ignore-integrated task perform well at the Store vs Ignore input gating tasks, but not the output 307 gating task, as expected, while models trained on the ignore-separated task perform poorly on both 308 input and output gating (3b). These results strongly support the intuition that while Transformers 309 have good inductive biases for learning role-addressable gating, it does not come "for free", and 310 emerges only when demanded by the training task – the same contrast in demands that demonstrated 311 advantages of frontostriatal gating (O'Reilly & Frank, 2006).

312 313

314

274 275

276

277 278

279

280

281

282

283

284

285

287

289

290

291

296

297

298

301

304

4.3 **EMERGENCE OF GATING PREDICTS TASK PERFORMANCE**

Is learning this gating policy useful for succeeding on the task? To answer this question, we first 315 compare models with the same hyperparameters across different random seeds. We train 20 new 316 models, each with a different random initialization, and measure both training loss and test set 317 accuracy. 5 of the models succeed 100% of the time, and the other 15 models succeed between 318 94%-99.99% of the time, with a mean of 97.72% and a standard deviation of 2.03. We measure the 319 intrinsic quality of the input and output gating subtasks using the path patching metrics described in 320 Section 3.3. 321

Figure 4 shows the 5 runs that reach 100% accuracy on the test data as well as 5 randomly selected 322 runs that do not. Two trends stand out. First, models which score a perfect test accuracy appear to 323 succeed at the gating subtasks more readily than models which do not. Of the former, 3 of 5 models



Figure 4: **Model Performance** over training on patching subtasks. Each graph contains an individual model's training loss (solid line) and subtask accuracy (dashed line, between 0 and 1) over time; the line's color corresponds to whether the model reaches 100% accuracy on the general test set.

338

339

340

reach 100% accuracy on both subtasks readily, plateauing less than halfway through training. In
contrast, models that make errors in the test set also do not reach such immediate success at the
subtasks (including the 10 not pictured in this graph); in fact, many categorically fail, scoring as low
as 49% accuracy. These results do not indicate that this class of models' representations are useless
for the task– they all score between 94% and 99.99%, well above chance performance.

The second trend is that many models across both classes have a sharp decline in training loss, which correlates with a similarly steep increase in accuracy on both subtasks. We interpret this phase transition as suddenly learning a gating mechanism. Models that do not exhibit phase transitions to the same degree take longer to fit the task, and do not reach high subtask accuracy.

352 While the above results suggest it is possible for a model to achieve fine accuracy in distribution 353 without learning a gating policy, we hypothesize that models which do learn a more general gating 354 strategy will better generalize to out of distribution examples. To assess this, we held out a subset 355 of symbols from each register (randomized which symbol would be held out from which register). 356 We tested these models on a challenge dataset which included examples of in-distribution pair-357 ings (register-symbol pairings that were seen in training) and out-of-distribution pairings (register-358 symbol pairings that were not seen during training). When holding out 5% of the symbols, the 359 models learn a robust input and output gating strategy, and notably, perform on average at 99.7% for 360 out-of-distribution symbol-register pairings. In contrast, the models trained on the split set controls do not learn robust gating strategies (despite performing perfectly at the trained task) and accord-361 ingly perform more poorly on the out of distribution examples: 77.2% (ignore integrated) and 88.6% 362 (ignore separated).¹ 363

364 Note that the Split Set (Ignore Separated) task depends least on learning an effective gating strategy. These models learn the task very quickly and show no ability to perform on the gating subtask, and instead likely learn heuristic solutions. Interestingly, when tested on the challenge dataset (which 366 includes examples of both in distribution and out of distribution register-symbol pairs), these mod-367 els show a large disruption in their ability to perform on in-distribution sequences (80.8%). This 368 insinuates that the strategy learned by these models is highly dependent on memorization and by 369 adding in new elements, the strategy fails. The out-of- distribution accuracy is slightly higher than 370 the in-distribution accuracy and future work could try to better understand the effects of perturbing 371 a brittle model. 372

373

¹One concern is that in the 5% held out, each register is trained on 49 symbols while in the split set (ignore integrated), each register is trained on 25 symbols. While all other parameters are held constant, this difference might be big enough to account for the large difference in generalization. To account for this, we ran another set of simulations holding out 60% of the symbols (each register is trained on 35 symbols). These models robustly learn a gating policy. Even with a drastic increase in the number of held out symbols from 5% to 60%, the out-of-distribution accuracy is still very high compared to either of the split set controls.

070			
370	Experiment	Out of Distribution	In Distribution
379	5% Held Out	99.7%	99.8%
380	60% Held Out	95.3%	99.3%
381	Split Set (Ignore Integrated)	77.2%	94.6%
382	Split Set (Ignore Separated)	88.6%	80.8%

Table 1: Accuracy for In Distribution and Out of Distribution Symbol-Register Pairings in Challenge Dataset Experiments and their associated performance on the challenge data set. The challenge dataset includes a mix of in-distribution and out-of-distribution symbol-register pairings. This table breaks down the accuracy based on if the symbol-register pairing was in the training (In-Distribution) or not seen during training (out-of-Distribution).

388 389 390

391 392

383

384

385

386

387

4.4 GATING POLICY TRANSFERS TO INCREASED TASK DEMANDS

Thus far, the experiments have focused on tasks with two registers, mimicking that used in the ref-393 erence back-2 task (Rac-Lubashevsky & Frank, 2021). However, human working memory has a 394 capacity of about 3-4 items (Cowan, 2008), albeit with vigorous debates questioning whether this 395 limit is discrete or continuous (Zhang & Luck, 2008; Wei et al., 2012; Luck & Vogel, 2013). More 396 recent models and data suggest a "chunking" hybrid between the two, whereby multiple memo-397 randa can compete for shared continuous resources within discrete slots (Nassar et al., 2018; Soni 398 & Frank, 2024). Chunking increases *effective* capacity, allowing more items to be remembered at 399 the cost of precision of some of the items. Notably, in frontostriatal gating models, when the num-400 ber of registers to manage was larger than two, networks given limited memory allocation but with 401 chunking capabilities performed better than those that were allocated as many PFC populations as items to store (Soni & Frank, 2024). The reason for this seeming paradox is *credit assignment*: as 402 the number of PFC populations (stripes) increases, the gating management problem becomes more 403 challenging - the network has to learn to route each item to distinct populations and to also learn to 404 read out from the corresponding population for a given probe. Moreover, these learning problems 405 are interdependent. 406

These limitations are computational rather than anatomical and stem from the learning process.
We hypothesize that they would also manifest in Transformers, even though Transformers have no
inherent memory demands (since all information is available in the context window). Specifically
because of the flexibility and expressivity of Transformers, we predicted that by increasing task
demand, the network would learn the task, but struggle to do so. We further predict that models that
learn mechanistic solutions will generalize better to new tasks.

To test this hypothesis, we increased the number of registers from 2 to 3. First, we confirmed that 413 asymptotic accuracy dropped to 95.4%. This result is qualitatively the same even when training 414 for twice the number of epochs. This is non-trivial given that we are just adding one register. We 415 predicted that the degree to which the transformer solves the task is related to heuristics and mem-416 orization rather than gating in these cases. We predicted that if we first pretrained the model to 417 effectively manage two registers, the network will be able to scaffold the learned gating strategy to 418 learn the three register task more robustly. We further predicted that this pretraining would only be 419 useful if pretraining encouraged gating strategies. Our results in Figure 5 support both of these con-420 clusions. Not only do pretrained networks on the original reference back-2 task exhibit higher gating 421 sub-task accuracy 5, but networks that were pretrained with the original two register task showed 422 very rapid learning in the three register task in the first few epochs. Moreover, this pretraining was 423 far less effective when it did not encourage gating (the split symbol control from above).

424 425

426 427

5 SUMMARY AND DISCUSSION

In this work, we investigate Transformer models for emergence of a learned *gating mechanism*; a network component performing role-addressable gating, similar to that in working memory of humans. We observe that the model learns a gating policy and find that task performance is correlated with gating ability. Our results show how learning gating mechanisms is one way Transformers can excel at tasks that require executive function.



445 Figure 5: 3 Register Task and Pretraining a)At the end of training on the 3 Register Task, models 446 with 2 Register Task pretraining perform the best on output gating tasks, insinuating learning of a mechanistic solution. b) Accuracy curves through training (after the pretraining) show a stark dif-447 ference between pretraining on the normal 2 register task and the control task. Models pretrained on 448 the 2 register task, generalize quickly (first few epochs) and with high accuracy (above 95%) on this 449 new task. At the end of training, no pretraining models are below 95% accuracy - which is below 450 the accuracy that 2 register task pretraining models showed at the beginning of training. Control 451 pretraining models show some generalization (due to learning of basic elements e.g. symbols) but 452 do not generalize to the same degree. Models pretrained with the 2 register tasks, which should 453 encourage a mechanistic solution, perform well at the new 3 register task, showing a superior ability 454 to generalize. 22 random seeds were run for each experimental condition and the results here an 455 average of those models. See Appendix A.5) for comparison on input and output gating task accu-456 racies.

- 457
- 458 459

The Transformer models are capable of making use of key composition for input gating and query 460 composition for output gating on the task. We find that making precise corruptions to specific ar-461 chitectural elements of the network causes the model's prediction to change from Same to Different 462 or vice versa, indicating that those components are causally responsible for the gating mechanism. 463 The architectural biases of attention within the vanilla Transformer model lend themselves well to 464 representing role-addressable content. The learnable nature of keys, queries, and values allows the 465 model to learn to create internal representations. These representations can be learned to signify roles and addresses, mimicking the variable binding and input / output gating mechanisms in bi-466 ological neural networks (O'Reilly & Frank, 2006; Frank & Badre, 2012; Collins & Frank, 2013; 467 Kriete et al., 2013). 468

469 When we trained more models on the task, we found that the models which perform best on the 470 task correlate with the markers of gating we observed in our circuit analysis, and that the learning trajectory shows a steep decrease in training loss and a steep rise in patched subtask accuracy 471 simultaneously, suggesting that the model has learned a gating policy at that time. Both findings 472 are analogous to those of Frank & Badre (2012), in which they find that networks which learned 473 a hierarchical gating policy performed better at a hierarchical learning task, and humans that learn 474 this policy also show a sharp decrease in loss when they discover it. There is still more work to 475 be done to better understand the models that don't learn the gating policy: what types of solutions 476 do they learn? How do we push models towards a mechanistic solution by hyperparameter tuning? 477 Understanding grokking phase planes in a similar manner as Liu et al. could be informative. 478

Nevertheless, we show that a critical factor controlling the learned gating policy is the task demands.
We found that experimental conditions in which biological networks exhibit gating advantages were similarly needed to give rise to learned gating policies in Transformers. Conversely, when the task places less demands on role-addressable gating (our split symbol control conditions), the models 1) did not learn gating policies, 2) were less able to generalize to out-of-distribution pairs, and 3) were less able to rapidly acquire tasks with higher gating demands (three registers).

485 We show that models that learn and solve the tasks using gating mechanisms are better at generalization. Further we show how learning brittle and memorized solutions causes some models to falter at even in-distribution pairings when they are mixed with out-of-distribution examples. These results
can be used to inform larger models that are needed for better and faster generalization. Future work
can characterize what kinds of mistakes are common within the mechanistic and heuristic models to
further characterize these models and promote better generalizability.

490 Ultimately, finding connections between emergent behavior of Transformer models and human 491 working memory serves to benefit both computational cognitive neuroscience and artificial intel-492 ligence. Although Transformer models themselves are limited in their biological plausibility, in this 493 setting they learned behavior mimicking the functionality of working memory, and their application 494 within computational models of the brain should be further explored. From the perspective of arti-495 ficial intelligence, understanding the strengths and limitations of Transformer models on executive 496 function tasks may inform model analysis across the many diverse settings in which these models are applied. 497

498 499 6 T

500

525

526

527

528

533

534

6 LIMITATIONS

501 While this paper draws parallels between the Transformer architecture and the brain, it is important 502 to emphasize that there are some significant differences between how Transformers and humans 503 solve tasks. In particular, because Transformers can attend to any part of the sequence when cre-504 ating a representation, they are not limited by memory constraints. Transformers can solve tasks 505 that would push the limits of human working memory (but it should be noted that Transformers require disproportionately large amounts of training data to do so). Nevertheless, we hypothesized 506 that Transformers might still learn effective gating policies that mimic those in frontostriatal net-507 works. Moreover, as briefly reviewed in the introduction, working memory capacity in humans and 508 biological computational models is not limited primarily by the memory demand per se, but rather 509 by the difficulty of the credit assignment problem for learning how to manage role-addressable stor-510 age and access of multiple items in memory. A fundamental bridge between Transformers and other 511 models of WM (and humans themselves) would be if Transformers also needed to overcome the 512 credit assignment learning problem, despite an unlimited memory capacity. 513

- 514 515 REFERENCES
- David Badre and Michael J Frank. Mechanisms of hierarchical reinforcement learning in cortico–
 striatal circuits 2: Evidence from fmri. *Cerebral cortex*, 22(3):527–536, 2012.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Bernhard Baier, Hans Otto Karnath, Marianne Dieterich, Frank Birklein, Carolin Heinze, and Notger G. Müller. Keeping memory clear and stable the contribution of human basal ganglia and
 prefrontal cortex to working memory. *Journal of Neuroscience*, 30, 2010. ISSN 15292401. doi:
 10.1523/JNEUROSCI.1513-10.2010.
 - Apoorva Bhandari and David Badre. Learning and transfer of working memory gating policies. Cognition, 172:89–100, 2018. ISSN 0010-0277. doi: https://doi.org/10.1016/j.cognition. 2017.12.001. URL https://www.sciencedirect.com/science/article/pii/s0010027717303037.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
 - Mikhail S Burtsev, Yuri Kuratov, Anton Peganov, and Grigory V Sapunov. Memory transformer. arXiv preprint arXiv:2006.11527, 2020.
- Cristian Buc Calderon, Tom Verguts, and Michael J Frank. Thunderstruck: The acdc model of flexible sequences and rhythms in recurrent neural circuits. *PLoS Computational Biology*, 18(2): e1009854, 2022.
- 539 Christopher H Chatham, Michael J Frank, and David Badre. Corticostriatal output gating during selection from working memory. *Neuron*, 81(4):930–942, 2014.

540 Anne GE Collins and Michael J Frank. Cognitive control over learning: creating, clustering, and 541 generalizing task-set structure. *Psychological review*, 120(1):190, 2013. 542 Nelson Cowan. Chapter 20 what are the differences between long-term, short-term, and working 543 memory?, 2008. ISSN 00796123. 544 Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdi-546 nov. Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint 547 arXiv:1901.02860, 2019. 548 549 Michael J Frank and David Badre. Mechanisms of hierarchical reinforcement learning in corticos-550 triatal circuits 1: computational analysis. Cerebral cortex, 22(3):509-526, 2012. 551 Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model 552 behavior with path patching. arXiv preprint arXiv:2304.05969, 2023. 553 554 S Hochreiter. Long short-term memory. Neural Computation MIT-Press, 1997. 555 556 Trenton Kriete, David C Noelle, Jonathan D Cohen, and Randall C O'Reilly. Indirection and symbol-like processing in the prefrontal cortex and basal ganglia. *Proceedings of the National* Academy of Sciences, 110(41):16390–16395, 2013. 558 559 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. nature, 521(7553):436-444, 560 2015. 561 562 Ziming Liu, Ouail Kitouni, Niklas Nolte, Eric J. Michaud, Max Tegmark, and Mike Williams. To-563 wards understanding grokking: Aneffective theory of representation learning. NeurIPs, 2022. Steven J. Luck and Edward K. Vogel. Visual working memory capacity: From psychophysics and 565 neurobiology to individual differences, 2013. ISSN 13646613. 566 567 Sam Whitman McGrath, Jacob Russin, Ellie Pavlick, and Roman Feiman. How can deep neural 568 networks inform theory in psychological science? Current Directions in Psychological Science, 569 pp. 09637214241268098, 2024. 570 Fiona McNab and Torkel Klingberg. Prefrontal cortex and basal ganglia control access to working 571 memory. Nature Neuroscience, 11, 2008. ISSN 10976256. doi: 10.1038/nn2024. 572 573 Neel Nanda and Joseph Bloom. Transformerlens. https://github.com/neelnanda-io/ 574 TransformerLens, 2022. 575 576 Matthew R. Nassar, Julie C. Helmers, and Michael J. Frank. Chunking as a rational strategy for lossy 577 data compression in visual working memory. Psychological Review, 125, 2018. ISSN 0033295X. doi: 10.1037/rev0000101. 578 579 Chris Olah. Mechanistic interpretability, variables, and the importance of interpretable bases. 580 https://www.transformer-circuits.pub/2022/mech-interp-essay, 2022. 581 582 Randall C O'Reilly and Michael J Frank. Making working memory work: a computational model 583 of learning in the prefrontal cortex and basal ganglia. Neural computation, 18(2):283–328, 2006. 584 Judea Pearl. Direct and indirect effects. In Proceedings of the Seventeenth Conference on Uncer-585 tainty in Artificial Intelligence, UAI'01, pp. 411-420, San Francisco, CA, USA, 2001. Morgan 586 Kaufmann Publishers Inc. ISBN 1558608001. 588 Rachel Rac-Lubashevsky and Michael J Frank. Analogous computations in working memory in-589 put, output and motor gating: Electrophysiological and computational modeling evidence. PLoS computational biology, 17(6):e1008971, 2021. Rachel Rac-Lubashevsky and Yoav Kessler. Dissociating working memory updating and automatic 592 updating: The reference-back paradigm. Journal of Experimental Psychology: Learning, Mem-593 ory, and Cognition, 42(6):951, 2016.

- A. Soni and M. J. Frank. Adaptive chunking improves effective working memory capacity in a prefrontal cortex and basal ganglia circuit. *eLife*, 13:RP97894, 2024. URL https: //elifesciences.org/reviewed-preprints/97894.
- Michael T. Todd, Yael Niv, and Jonathan D. Cohen. Learning to use working memory in partially observable environments through dopaminergic reinforcement. In *Advances in Neural Information Processing Systems 21 - Proceedings of the 2008 Conference*, 2009.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- Edward K. Vogel, Andrew W. McCollough, and Maro G. Machizawa. Neural measures reveal individual differences in controlling access to working memory. *Nature*, 438, 2005. ISSN 14764687. doi: 10.1038/nature04171.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*, 2022.
 - Yau-Shian Wang, Hung-Yi Lee, and Yun-Nung Chen. Tree transformer: Integrating tree structures into self-attention. *arXiv preprint arXiv:1909.06639*, 2019.
- Ziqiang Wei, Xiao Jing Wang, and Da Hui Wang. From distributed resources to limited slots in
 multiple-item working memory: A spiking network model with normalization. *Journal of Neuroscience*, 32, 2012. ISSN 02706474. doi: 10.1523/JNEUROSCI.0735-12.2012.
 - Weiwei Zhang and Steven J. Luck. Discrete fixed-resolution representations in visual working memory. *Nature*, 453, 2008. ISSN 14764687. doi: 10.1038/nature06860.
- 622 623 624

627

620

621

614

615

616

A APPENDIX

A.1 SELF-ATTENTION IN TRANSFORMERS

Transformers are powerful language models which create contextualized representations of the input sequences. They learn to predict the next token one at a time using an "attention" mechanism that scans other tokens in the sequence for relevant information (Bahdanau et al., 2014). A common practice (that is used here) is to mask any future tokens and so predictions and representations must be made by the current token and any previously seen tokens. These models are able to learn and represent complex sequence modelling tasks.

For a given prediction, Transformer attention generates three separate vectors for each token in 634 the sequence: a query, key, and value (q, k, v). The key vector is a set of tokens that the model 635 has learned are most relevant to the token at hand. The query vector scans the tokens in the key 636 vectors and calculates how much the current prediction should "attend" to each of those tokens. 637 Then, the value vectors at those positions are multiplied by the corresponding weights, summed 638 up, and added to the next representation: for token i at layer j, the contextual representation is 639 $\sum_{k} q_i^j \cdot k_k^j * v_k^j$. Thus, the next token prediction includes earlier sequential information by combining 640 the value vectors from previous tokens. In other words, Transformer attention can be viewed as a 641 read/write mechanism: for a given token, the queries and keys dictate which tokens to read from, the 642 values are the content that is read (proportional to the attention calculated by the keys and queries), 643 and the summed content is written to a new representation at the given token. As we shall see below, 644 the comparison to role-addressable input and output gating operations is evident. The key vectors 645 form addresses analogous to the PFC "stripes" (neural populations that support variable binding in memory). The learned key construction determines the address to store an item and is thus analogous 646 to input gating. Conversely, the query vectors determine which addresses are accessed, and are thus 647 analogous to output gating.

648 A.2 **TEXTUAL REFERENCE-BACK-2 TASK** 649

650 The textual reference-back task requires making same/different judgments between incoming sym-651 bols assigned to a particular "register" in memory, with respect to those seen previously and linked to those same registers. Like the original tasks, the textual reference-back task is sequential, and 652 requires independent updating and maintenance of two memory registers, each containing one of 653 S arbitrary symbols at a time. At the beginning of each sequence, each register is initialized indi-654 vidually to one $s \in S$ (the pool of symbols is shared between registers, which was shown to more 655 substantively tax gating mechanisms in (?).) Each sequence is composed of L tuples, each contain-656 ing register address Reg_i, symbol Sym_i, same/different label Ans_i, and update instruction Ins_i. For 657 a tuple $i \in L$, the **answer** Ans_i is a binary value that is either Same if symbol Sym_i is currently 658 stored in the register with address Reg_i , or Different otherwise. The **update instruction** Ins_i also 659 takes one of two values (ignore or store), evenly distributed. If the instruction is ignore, then 660 the model still needs to make the same/different determination with respect to the stored reference, 661 but the new symbol should not update the register content (i.e., the reference remains unperturbed). 662 If it is store, then from that point on in the sequence, Sym_i is stored in the register with address Reg, until otherwise updated. An example is shown in Fig. 1. 663

664 We implement each reference-back task example in our data as a single sequence, and measure 665 models' ability to predict Same versus Different for each Ans_i. Each sequence has 10 same/different 666 answers, and we generate 100,000 train, 1,000 validation, and 1,000 held-out test sequences.

667 The class balance of same to different answer labels in the train/test datasets is roughly 1:2, 668 making a "maximum class" heuristic solution 0.66 accuracy, 0.33 precision, and 0.5 recall. We test 669 several other heuristics, the strongest of which is predicting same if another tuple including Store 670 and the target register and target symbol exists in the sequence, which scores 0.80 accuracy, 0.82 671 precision, and 0.85 recall.

672

674

673 A.3 MODELS AND TRAINING

675 We train small, attention-only Transformer models from scratch on our task. Our models contain 676 two decoder-only layers, each with two heads, and no multilayer perceptrons or layer normalization, 677 followed by a linear "unembed" layer to project the output of the last decoder into the space of the 678 entire vocabulary at each timestep In practice, only 'same' and 'different' are ever predicted. Our network uses absolute positional embeddings (Vaswani et al., 2017). The vocabulary contains all 679 possible tokens, represented individually with embedding size E. Models are trained to predict the 680 next token with the language modelling objective: if the model is predicting Ans_c, it will have access 681 to all $(Ins_i, Reg_i, Sym_i, Ans_i)$ tuples where i < c, as well as Ins_c, Reg_c , and Sym_c . However, the 682 models only receive loss at positions where a same/different token must be predicted (this is similar 683 to the reward function applied in frontostriatal gating networks; (O'Reilly & Frank, 2006; Soni & 684 Frank, 2024)). Furthermore, each layer gets a causal attention mask– when constructing each token 685 representation, it cannot look ahead at tokens further down the sequence. 686

The models are trained over 60 epochs of the 100k training data points, learning from 6 million 687 examples in total. Models are evaluated on their accuracy (whether the correct Ans_i is predicted for 688 each tuple i), measured in precision and recall, as well as the same versus different token logit 689 difference. 690

691

A.4 PATH PATCHING 692

693 Path-patching involves making a incisive edit to the representations of a trained model and ob-694 serving how the model's behavior is affected, allowing one to infer the computations implemented 695 within individual attentional heads (see Fig. 6). Path-patching requires a minimal pair of examples: 696 the "clean" example and the "corrupted" example, in which one token from the clean example is 697 changed, as well as the correct label. Given representations from the model for both the clean and 698 the corrupted examples (the blue and orange components in the figure), we can chose a specific 699 component anywhere in the model (referred to as the "sender"), and replace the clean representation with the corrupted one at that specific component. These embeddings will be received by the next 700 layer ("receiver"), thereby "patching" the path. From there, the patched model will compute the new 701 prediction.



Figure 6: **Graphical Path Patching** Graphical diagram of the path-patching process. Attention heads are represented as circles (layer,head index), and contextual representations of each token (as well as the next token prediction) are represented as rectangles.

746

720

721

In the figure, we send from layer 0, head 0 and 1 to layer 1, both heads 0 and 1. All clean representations that are not along this path are not modified and are unaffected by the patch. The model then recomputes all representations after the receiver (the "patched" representations), and arrives at a new prediction. If the model output matches the corrupt prediction rather than the clean one, that prediction is causally dependent on the path from sender to receiver. See (Wang et al., 2022) and (Goldowsky-Dill et al., 2023) for a more comprehensive review of path-patching methods.

We perform a small hyperparameter search and select a model that reaches 100% accuracy on the
held-out test data for further analysis. We determine the circuit that the model uses through an array
of path-patching experiments with a simple minimal pairs paradigm. Our "sender" within pathpatching is always both attention heads at layer 0, and our "receiver" is always both attention heads
at layer 1.

- ⁷³⁵ We first establish that a Transformer model is able to succeed on the reference-back-2 task.
- At layer 0, the model learns to condense the task-critical information from each tuple into one embedding, at the position for Sym_i^2 . At this layer, the model pays 85.8% of total attention to the task-critical information to that tuple, and just 14.2% of attention to other tuples.

At layer 1, the attention heads learn to attend to the Sym_i key vector representing the tuple where information was last stored in the target register. The heads pay 70.2% of total attention to this tuple (the "stored" tuple), and only 29.8% of attention to all other tokens. This behavior is tied to the target register matching the register in the stored tuple, which is analogous to gating of the relevant role-addressable PFC stripe (O'Reilly & Frank, 2006; Soni & Frank, 2024). We focus our analysis on the Layer 1 representations which exhibit this learned gating policy, shown in Fig. 2.

A.5 INPUT VS. OUPUT GATING ACCURACIES, PRETRAINING

²Redundantly, the model does the same at the position for Ans_i . Through additional experimentation, we determine that this is a quirk of Transformer learning, and does not impact our analysis.



Figure 7: **Input and Output Gating Task Accuracy** Pretraining on the 2 register task leads to the highest accuracy on input and output gating subtasks - insinuating that these models learn the most mechansitic solutions. In general the input gating accuracy is lower, indidcating a differential role for each gating. There is human experimental evidence to suggest that output gating is harder to learn and is more crucial for better performance