

Discovering Interpretable Algorithms by Decompiling Transformers to RASP

Xinting Huang*



Aleksandra Bakalova*



Satwik Bhattamishra



William Merrill



Michael Hahn



What algorithms do Transformers internally implement?

Weiss, Gail, Yoav Goldberg, and Eran Yahav. "Thinking like transformers." International Conference on Machine Learning. PMLR, 2021.

Zhou, H., Bradley, A., Littwin, E., Razin, N., Saremi, O., Susskind, J. M., Bengio, S., and Nakkiran, P. What Algorithms can Transformers Learn? A Study in Length Generalization. In International Conference on Learning Representations, 2024.

Huang, X., Yang, A., Bhattamishra, S., Sarrof, Y., Krebs, A., Zhou, H., ... & Hahn, M. (2025, May). A formal framework for understanding length generalization in transformers. In International Conference on Learning Representations (Vol. 2025, pp. 58095-58179).

What algorithms do Transformers internally implement?

RASP_(Weiss et al, 2021) - a symbolic language that mimics computations in transformers.

```
19 select_last = select(indices, length-1, ==);
20 end_0 = aggregate(select_last,
21                 bal1==0 and bal2==0);
22
23 shuffle_dyck2 = end_0 and not had_neg;
```

Weiss, Gail, Yoav Goldberg, and Eran Yahav. "Thinking like transformers." International Conference on Machine Learning. PMLR, 2021.

Zhou, H., Bradley, A., Littwin, E., Razin, N., Saremi, O., Susskind, J. M., Bengio, S., and Nakkiran, P. What Algorithms can Transformers Learn? A Study in Length Generalization. In International Conference on Learning Representations, 2024.

Huang, X., Yang, A., Bhattamishra, S., Sarrof, Y., Krebs, A., Zhou, H., ... & Hahn, M. (2025, May). A formal framework for understanding length generalization in transformers. In International Conference on Learning Representations (Vol. 2025, pp. 58095-58179).

What algorithms do Transformers internally implement?

RASP_(Weiss et al, 2021) - a symbolic language that mimics computations in transformers.

```
19 select_last = select(indices, length-1, ==);
20 end_0 = aggregate(select_last,
21                 ball==0 and bal2==0);
22
23 shuffle_dyck2 = end_0 and not had_neg;
```

Transformers length-generalize if trained on tasks which are expressible in **(C-)RASP**_(Zhou et al., 2024; Huang et al., 2025).

Weiss, Gail, Yoav Goldberg, and Eran Yahav. "Thinking like transformers." International Conference on Machine Learning. PMLR, 2021.

Zhou, H., Bradley, A., Littwin, E., Razin, N., Saremi, O., Susskind, J. M., Bengio, S., and Nakkiran, P. What Algorithms can Transformers Learn? A Study in Length Generalization. In International Conference on Learning Representations, 2024.

Huang, X., Yang, A., Bhattamishra, S., Sarrof, Y., Krebs, A., Zhou, H., ... & Hahn, M. (2025, May). A formal framework for understanding length generalization in transformers. In International Conference on Learning Representations (Vol. 2025, pp. 58095-58179).

What algorithms do Transformers internally implement?

RASP_(Weiss et al, 2021) - a symbolic language that mimics computations in transformers.

```
19  select_last = select(indices, length-1, ==);
20  end_0 = aggregate(select_last,
21                ball==0 and bal2==0);
22
23  shuffle_dyck2 = end_0 and not had_neg;
```

Transformers length-generalize if trained on tasks which are expressible in (C-)RASP_(Zhou et al., 2024; Huang et al., 2025).

Length-generalizable transformers might internally implement RASP programs, but do they?

Weiss, Gail, Yoav Goldberg, and Eran Yahav. "Thinking like transformers." International Conference on Machine Learning. PMLR, 2021.

Zhou, H., Bradley, A., Littwin, E., Razin, N., Saremi, O., Susskind, J. M., Bengio, S., and Nakkiran, P. What Algorithms can Transformers Learn? A Study in Length Generalization. In International Conference on Learning Representations, 2024.

Huang, X., Yang, A., Bhattamishra, S., Sarrof, Y., Krebs, A., Zhou, H., ... & Hahn, M. (2025, May). A formal framework for understanding length generalization in transformers. In International Conference on Learning Representations (Vol. 2025, pp. 58095-58179).

Contributions

1. Method for automatically extracting short and interpretable RASP-like programs from trained Transformers.

To the best of our knowledge, this is the first method to convert naturally trained transformers to RASP.

Contributions

1. Method for automatically extracting short and interpretable RASP-like programs from trained Transformers.

To the best of our knowledge, this is the first method to convert naturally trained transformers to RASP.

2. Most direct evidence so far that transformers often internally implement interpretable RASP programs.

Decompilation Method

Main idea: **reparametrize and simplify.**

Programs are defined in D-RASP (Decompiled RASP).

D-RASP

1 $\alpha = \text{select}(k=v_k, q=v_q, \text{op}=\mathbf{A})$

$$\alpha(i, s) = v_q(i)^\top \mathbf{A} v_k(s)$$

2 $v = \text{aggregate}(s=\alpha_1 + \dots + \alpha_p, v=w)$

$$v(i) = \sum_{j \leq i} a_{i,j} w(j)$$

$$a_{i,s} = \frac{\exp(\alpha(i, s)_1 + \dots + \alpha(i, s)_p)}{\sum_{s' \leq i} \exp(\alpha(i, s')_1 + \dots + \alpha(i, s')_p)}$$

3 $v = \text{element_wise_op}(v_1, \dots, v_s, \text{func}=f)$

$$v(i) = f(v_1(i), \dots, v_s(i))$$

4 $p = \text{project}(\text{inp}=v, \text{op}=\mathbf{A})$

$$p(j) = \mathbf{A} \cdot v(j)$$

Decompilation Method

Step 1: Reparametrize

Theorem (informal). Given any GPT-2-style transformer satisfying Linear Layer Norm Assumption, we can rewrite it as a D-RASP program defining the same input-output map.

Decompilation Method

Step 1: Reparametrize

Neural network:

1 layer 1 head
transformer

D-RASP:

Program for 1-layer 1-head model

```

$$\alpha_{1,h} = \text{select}(k = \text{pos}, q = \text{token}, op = \mathbf{E}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{P})$$

$$\alpha'_{1,h} = \text{select}(k = \text{token}, q = \text{token}, op = \mathbf{E}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{E})$$

$$\alpha''_{1,h} = \text{select}(k = \text{pos}, q = \text{pos}, op = \mathbf{P}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{P})$$

$$\alpha'''_{1,h} = \text{select}(k = \text{token}, q = \text{pos}, op = \mathbf{P}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{E})$$

$$v_{(\text{pos},(1,h))} = \text{aggregate}(\alpha_{1,h} + \alpha'_{1,h} + \alpha''_{1,h} + \alpha'''_{1,h}, \text{pos})$$

$$v_{(\text{tok},(1,h))} = \text{aggregate}(\alpha_{1,h} + \alpha'_{1,h} + \alpha''_{1,h} + \alpha'''_{1,h}, \text{token})$$

$$v_{(\text{mlp}_1, \langle \rangle)} = \text{element-wise-op}(\text{token}, \text{pos}, v_{(\text{pos},(1,h))}, v_{(\text{tok},(1,h))}, \text{func} = f_{MLP,1})$$

$$p_{\text{token}} = \text{project}(\text{token}, op = \mathbf{U} \mathbf{E})$$

$$p_{\text{pos}} = \text{project}(\text{pos}, op = \mathbf{U} \mathbf{P})$$

$$p_{(\text{tok},(1,h))} = \text{project}(v_{(\text{tok},(1,h))}, op = \mathbf{U} \mathbf{V}_{1,h} \mathbf{E})$$

$$p_{(\text{pos},(1,h))} = \text{project}(v_{(\text{pos},(1,h))}, op = \mathbf{U} \mathbf{V}_{1,h} \mathbf{P})$$

$$p_{(\text{mlp}_1, \langle \rangle)} = \text{project}(v_{(\text{mlp}_1, \langle \rangle)}, op = \mathbf{U})$$

$$\text{prediction} = \text{softmax}(p_{\text{token}} + p_{\text{pos}} + p_{(\text{tok},(1,h))} + p_{(\text{pos},(1,h))} + p_{(\text{mlp}_1, \langle \rangle)})$$

```

Decompilation Method

Step 1: Reparametrize

Neural network:

1 layer 1 head
transformer

Program length is exponential in the
number of model's layers

D-RASP:

Program for 1-layer 1-head model

```

$$\begin{aligned} \alpha_{1,h} &= \text{select}(k = \text{pos}, q = \text{token}, op = \mathbf{E}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{P}) \\ \alpha'_{1,h} &= \text{select}(k = \text{token}, q = \text{token}, op = \mathbf{E}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{E}) \\ \alpha''_{1,h} &= \text{select}(k = \text{pos}, q = \text{pos}, op = \mathbf{P}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{P}) \\ \alpha'''_{1,h} &= \text{select}(k = \text{token}, q = \text{pos}, op = \mathbf{P}^T \mathbf{Q}_{1,h}^T \mathbf{K}_{1,h} \mathbf{E}) \\ v_{(\text{pos},(1,h))} &= \text{aggregate}(\alpha_{1,h} + \alpha'_{1,h} + \alpha''_{1,h} + \alpha'''_{1,h}, \text{pos}) \\ v_{(\text{tok},(1,h))} &= \text{aggregate}(\alpha_{1,h} + \alpha'_{1,h} + \alpha''_{1,h} + \alpha'''_{1,h}, \text{token}) \\ v_{(\text{mlp}_1, \diamond)} &= \text{element-wise-op}(\text{token}, \text{pos}, v_{(\text{pos},(1,h))}, v_{(\text{tok},(1,h))}, \text{func} = f_{MLP,1}) \\ p_{\text{token}} &= \text{project}(\text{token}, op = \mathbf{U}\mathbf{E}) \\ p_{\text{pos}} &= \text{project}(\text{pos}, op = \mathbf{U}\mathbf{P}) \\ p_{(\text{tok},(1,h))} &= \text{project}(v_{(\text{tok},(1,h))}, op = \mathbf{U}\mathbf{V}_{1,h} \mathbf{E}) \\ p_{(\text{pos},(1,h))} &= \text{project}(v_{(\text{pos},(1,h))}, op = \mathbf{U}\mathbf{V}_{1,h} \mathbf{P}) \\ p_{(\text{mlp}_1, \diamond)} &= \text{project}(v_{(\text{mlp}_1, \diamond)}, op = \mathbf{U}) \\ \text{prediction} &= \text{softmax}(p_{\text{token}} + p_{\text{pos}} + p_{(\text{tok},(1,h))} + p_{(\text{pos},(1,h))} + p_{(\text{mlp}_1, \diamond)}) \end{aligned}$$

```

Decompilation Method

Step 2.1: Simplify (pruning)

Neural network:

1 layer 1 head
transformer

D-RASP:

Program for 1-layer 1-head model

```
 $\alpha_{1,h} = \text{select}(k = \text{pos}, q = \text{token}, op = E^T Q_{1,h}^T K_{1,h} P)$   
 $\alpha'_{1,h} = \text{select}(k = \text{token}, q = \text{token}, op = E^T Q_{1,h}^T K_{1,h} E)$   
 $\alpha''_{1,h} = \text{select}(k = \text{pos}, q = \text{pos}, op = P^T Q_{1,h}^T K_{1,h} P)$   
 $\alpha'''_{1,h} = \text{select}(k = \text{token}, q = \text{pos}, op = P^T Q_{1,h}^T K_{1,h} E)$   
 $v_{(\text{pos}, \langle 1, h \rangle)} = \text{aggregate}(\alpha_{1,h} + \alpha'_{1,h} + \alpha''_{1,h} + \alpha'''_{1,h}, \text{pos})$   
 $v_{(\text{tok}, \langle 1, h \rangle)} = \text{aggregate}(\alpha_{1,h} + \alpha'_{1,h} + \alpha''_{1,h} + \alpha'''_{1,h}, \text{token})$   
 $v_{(\text{mp}_1, \langle \rangle)} = \text{element-wise-op}(\text{token}, \text{pos}, v_{(\text{pos}, \langle 1, h \rangle)}, v_{(\text{tok}, \langle 1, h \rangle)}, \text{func } f_{MLT,1})$   
 $p_{\text{token}} = \text{project}(\text{token}, op = UE)$   
 $p_{\text{pos}} = \text{project}(\text{pos}, op = UP)$   
 $p_{(\text{tok}, \langle 1, h \rangle)} = \text{project}(v_{(\text{tok}, \langle 1, h \rangle)}, op = UV_{1,h} E)$   
 $p_{(\text{pos}, \langle 1, h \rangle)} = \text{project}(v_{(\text{pos}, \langle 1, h \rangle)}, op = UV_{1,h} P)$   
 $p_{(\text{mp}_1, \langle \rangle)} = \text{project}(v_{(\text{mp}_1, \langle \rangle)}, op = UL)$   
 $\text{prediction} = \text{softmax}(p_{\text{token}} + p_{\text{pos}} + p_{(\text{tok}, \langle 1, h \rangle)} + p_{(\text{pos}, \langle 1, h \rangle)} + p_{(\text{mp}_1, \langle \rangle)})$ 
```

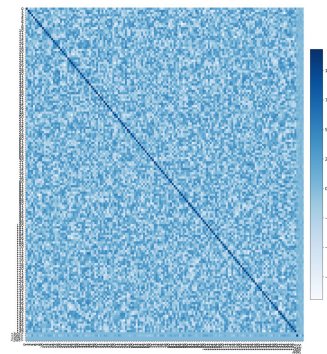
Decompilation Method

Step 2.2: Simplify (primitive matching)

...

$$a = \text{select}(q=\text{token}, k=v_{\langle tok, \langle (1,1) \rangle \rangle}, op=\mathbf{E}^T \mathbf{Q}_{2,1}^T \mathbf{K}_{2,1} \mathbf{V}_{1,1} \mathbf{E})$$

...



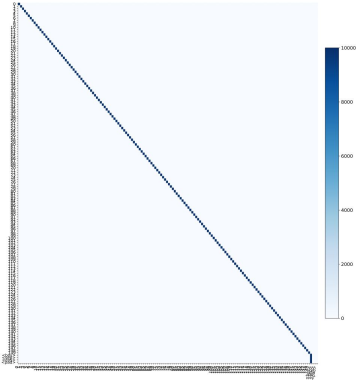
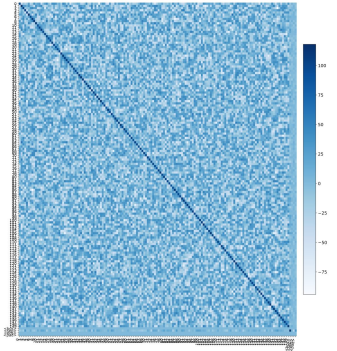
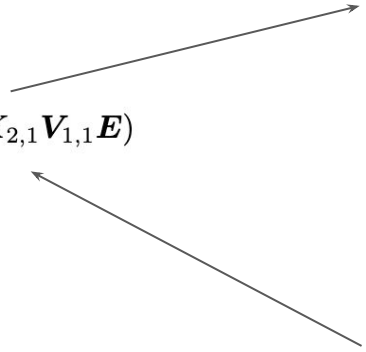
Decompilation Method

Step 2.2: Simplify (primitive matching)

...

```
a = select (q=token, k=v_{tok,((1,1))}, op= $\mathbf{E}^T \mathbf{Q}_{2,1}^T \mathbf{K}_{2,1} \mathbf{V}_{1,1} \mathbf{E}$ )
```

...



(b) Line 3: $op=(k==q)$, $special.op=(k==BOS)$

Decompilation Method

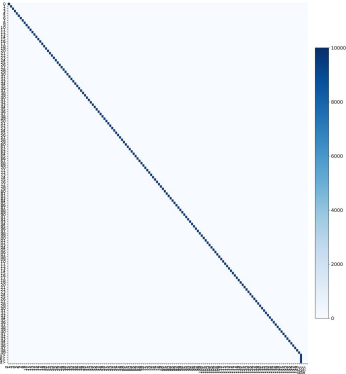
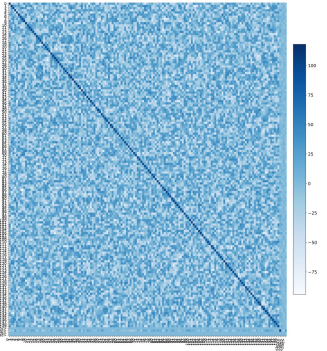
Step 2.2: Simplify (primitive matching)

...

```
a = select (q=token, k=v_{tok,((1,1))}, op= $\mathbf{E}^T \mathbf{Q}_{2,1}^T \mathbf{K}_{2,1} \mathbf{V}_{1,1} \mathbf{E}$ )
```

...

Replace



(b) Line 3: op=(k==q), special_op=(k==BOS)

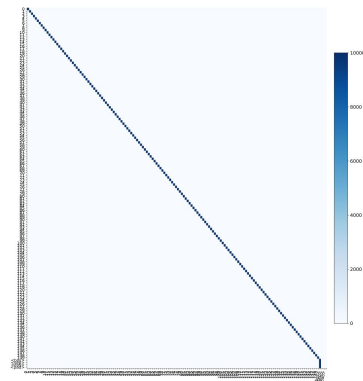
Decompilation Method

Step 2.2: Simplify (primitive matching)

...

$a = \text{select}(q=\text{token}, k=v_{\langle \text{tok}, \langle (1,1) \rangle \rangle}, \text{op}=\mathbf{E}^T \mathbf{Q}_{2,1}^T \mathbf{K}_{2,1} \mathbf{V}_{1,1} \mathbf{E})$

...



(b) Line 3: $\text{op}=(k==q)$, $\text{special_op}=(k==\text{BOS})$

Example

Decompiled program for MOST FREQUENT

A program extracted
from a 1-layer 1-head
transformer trained
on Most Frequent
task

Decompiled program for MOST FREQUENT

<bos> o b r o <sep>

Decompiled program for MOST FREQUENT

<bos> o b r o <sep>

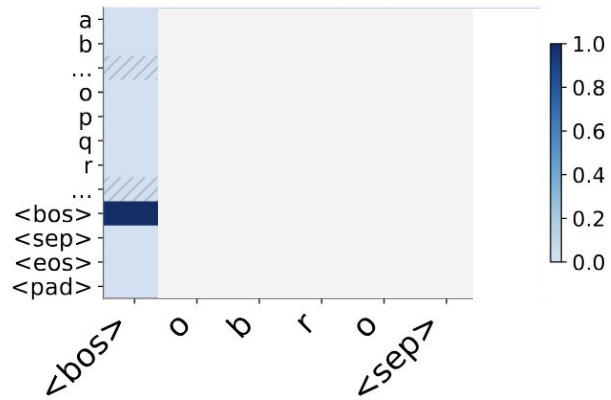
Decompiled program for MOST FREQUENT

<bos> o b r o <sep> o

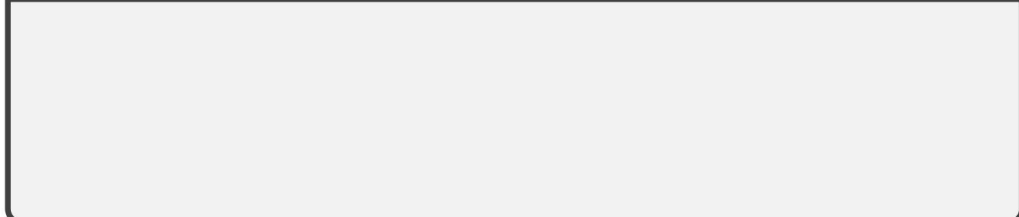
Decompiled program for MOST FREQUENT



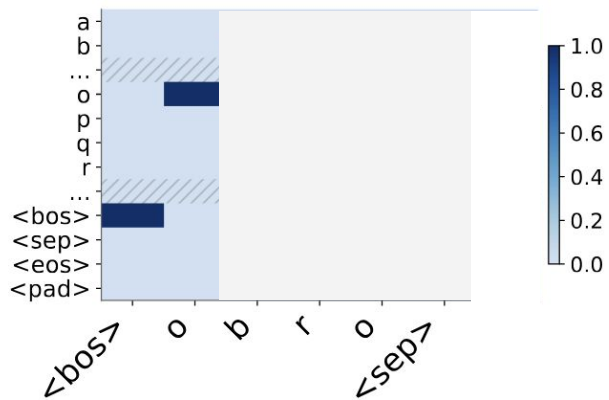
token



Decompiled program for MOST FREQUENT



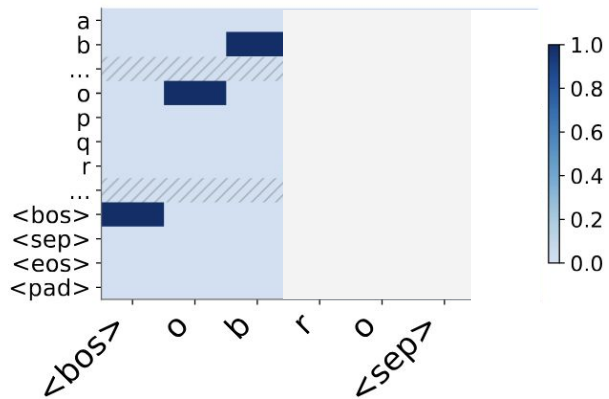
token



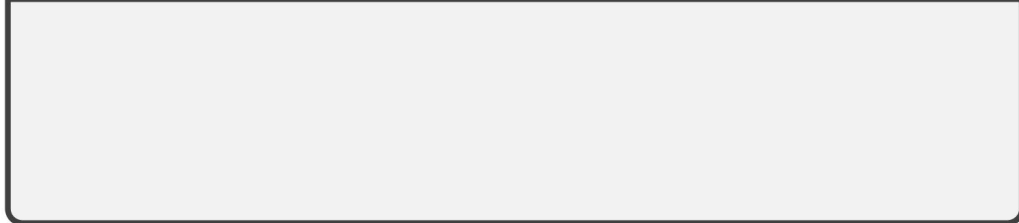
Decompiled program for MOST FREQUENT



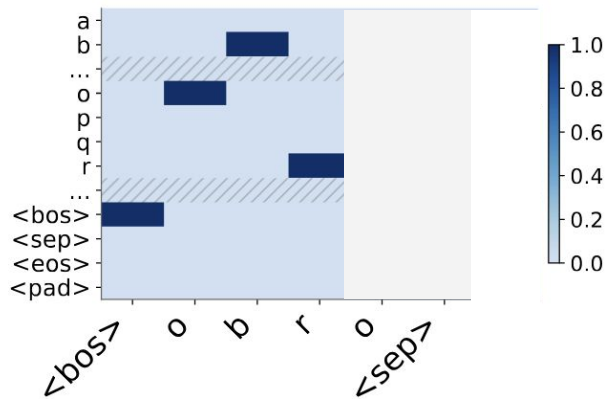
token



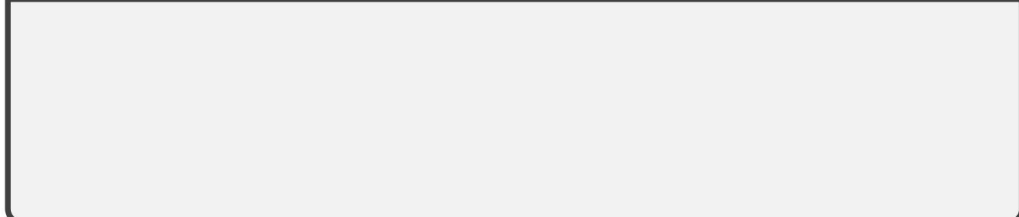
Decompiled program for MOST FREQUENT



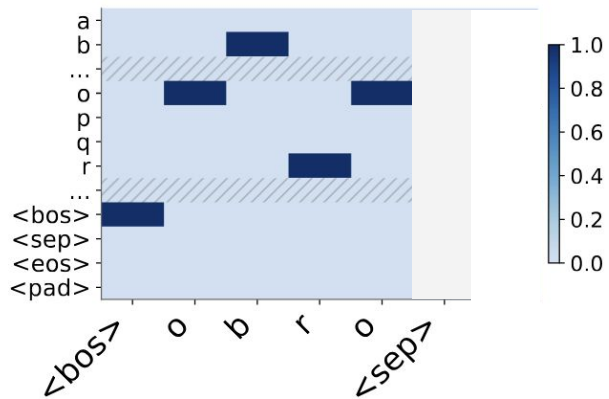
token



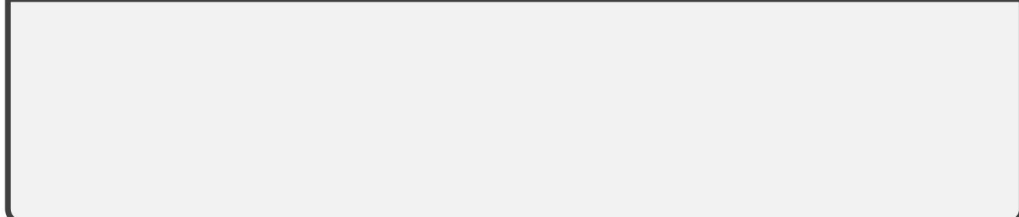
Decompiled program for MOST FREQUENT



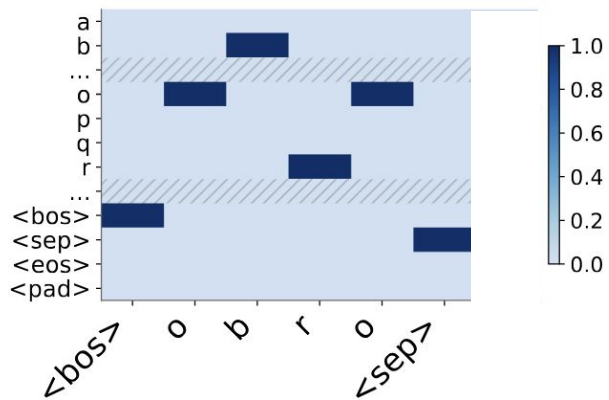
token



Decompiled program for MOST FREQUENT

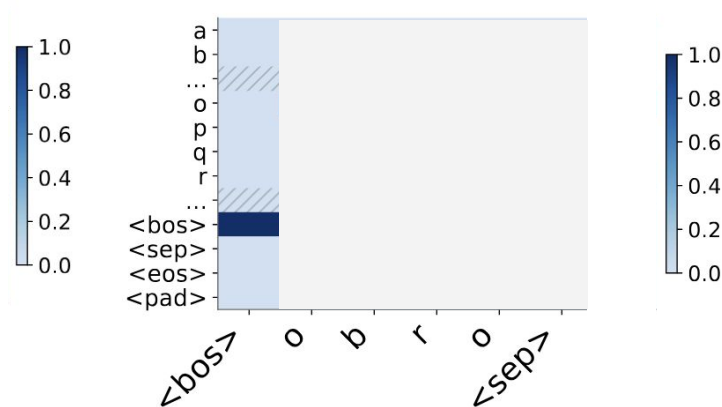
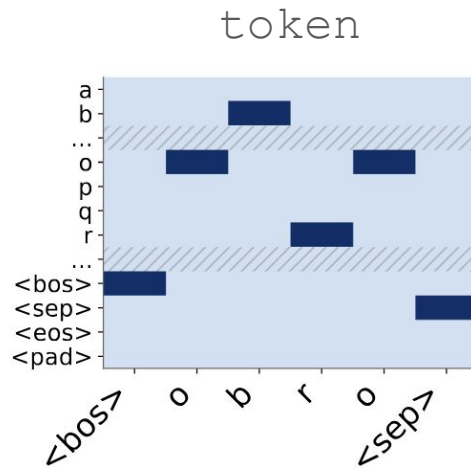


token



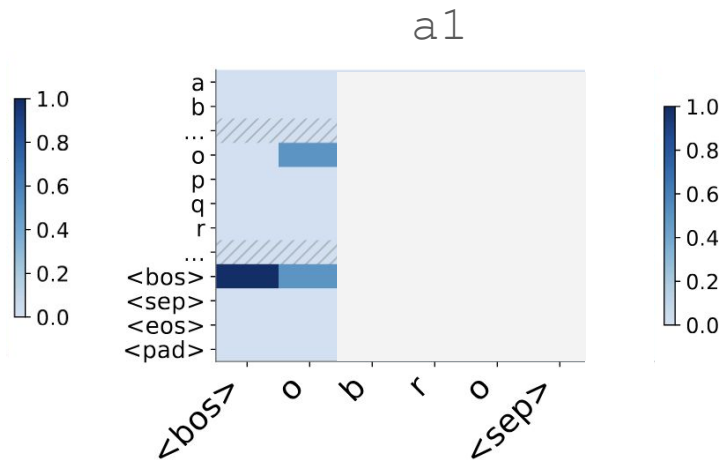
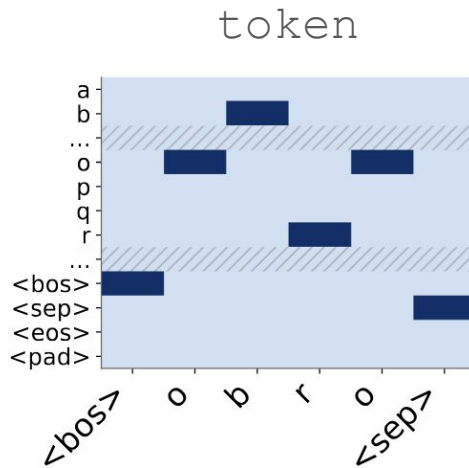
Decompiled program for MOST FREQUENT

```
1. a1 = aggregate(s=[], v=token)
```



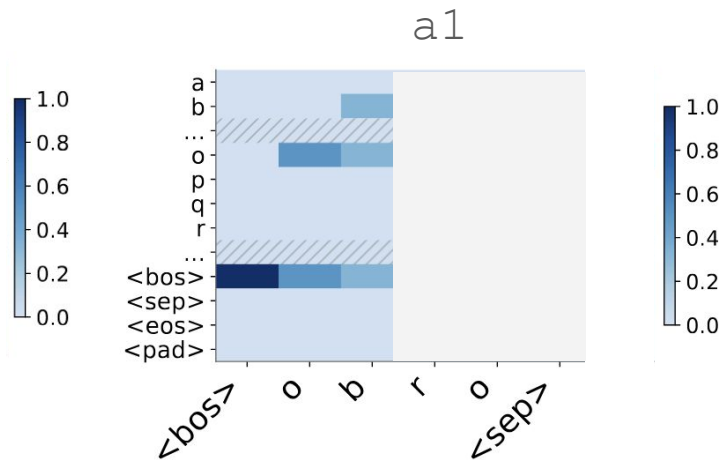
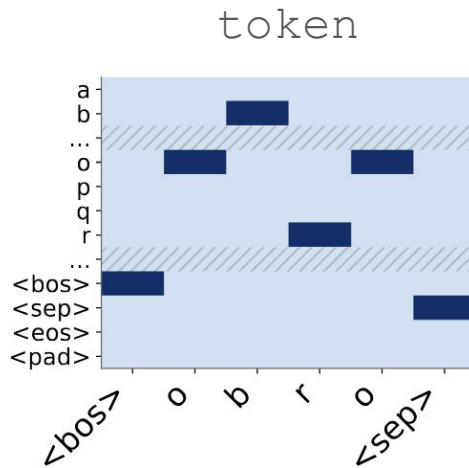
Decompiled program for MOST FREQUENT

```
1. a1 = aggregate(s=[], v=token)
```



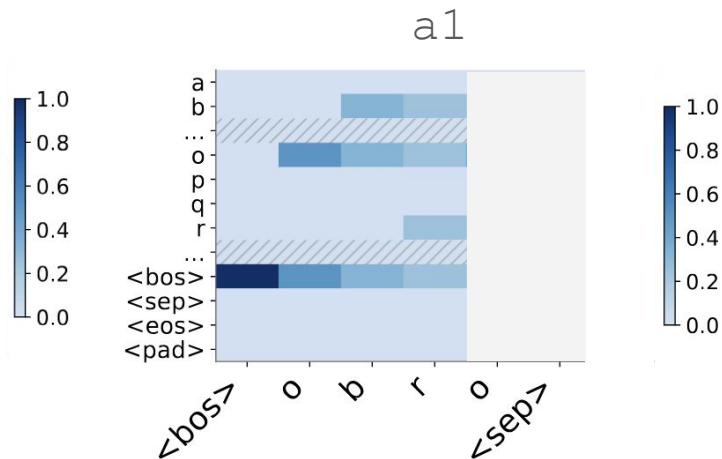
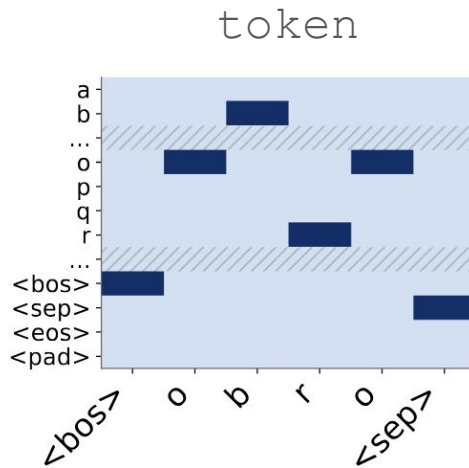
Decompiled program for MOST FREQUENT

```
1. a1 = aggregate(s=[], v=token)
```



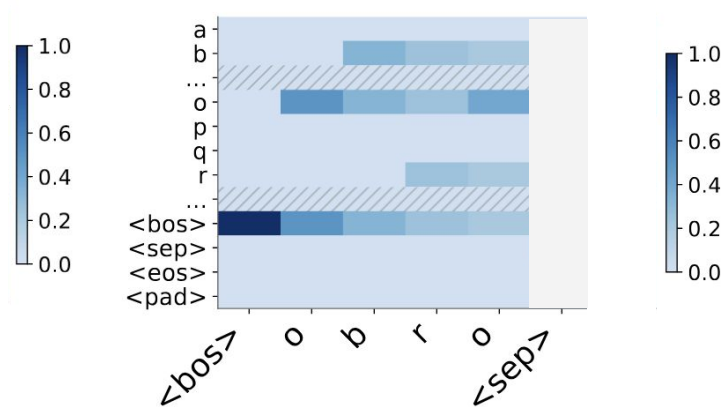
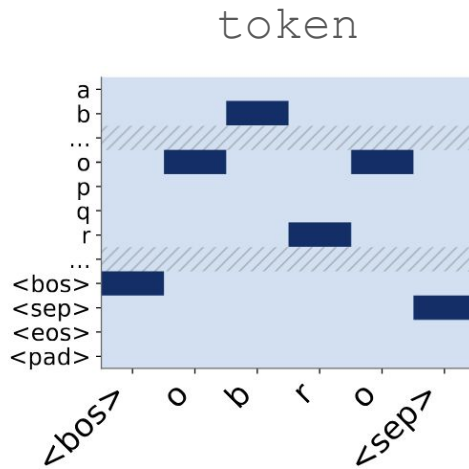
Decompiled program for MOST FREQUENT

```
1. a1 = aggregate(s=[], v=token)
```



Decompiled program for MOST FREQUENT

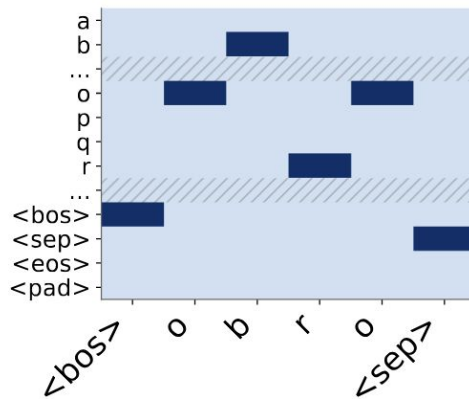
```
1. a1 = aggregate(s=[], v=token)
```



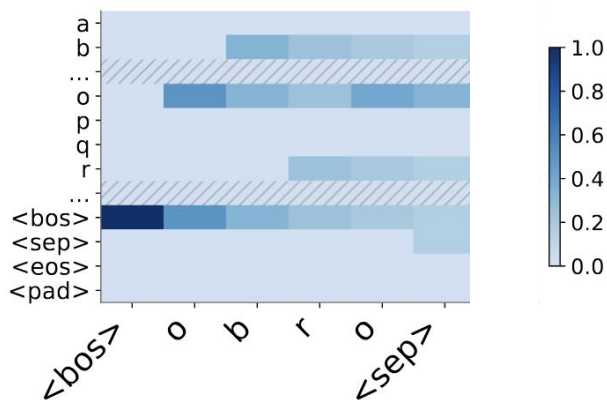
Decompiled program for MOST FREQUENT

```
1. a1 = aggregate(s=[], v=token)
```

token

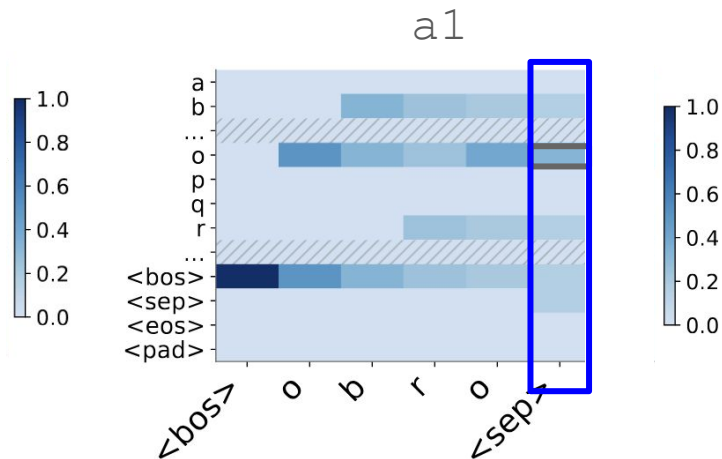
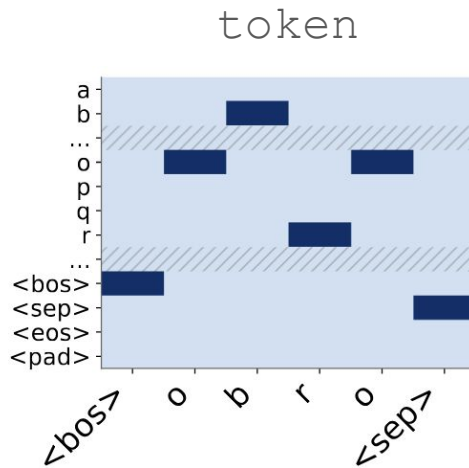


a1



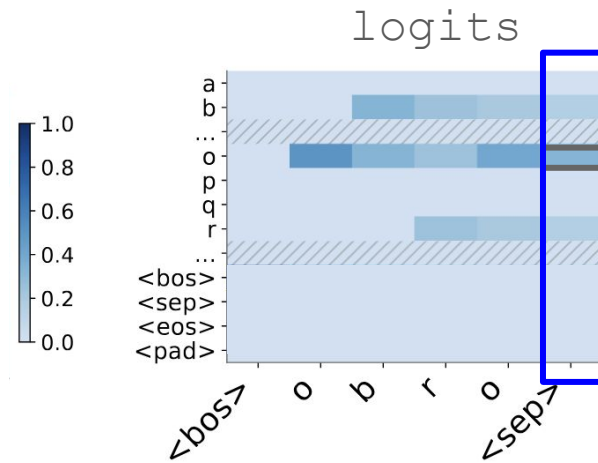
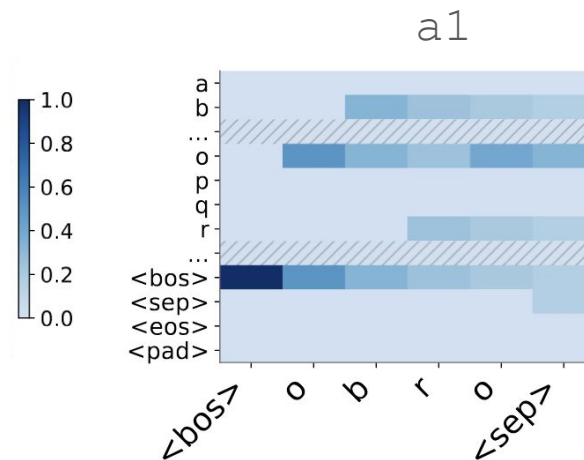
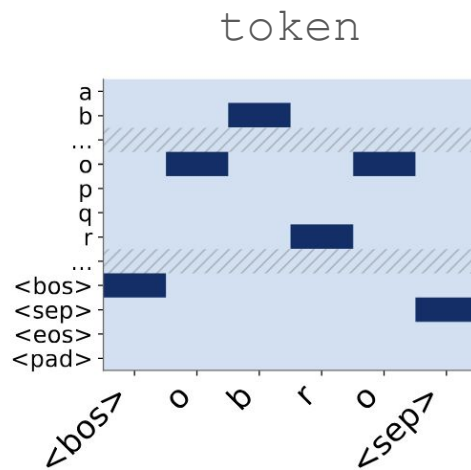
Decompiled program for MOST FREQUENT

```
1. a1 = aggregate(s=[], v=token)
```



Decompiled program for MOST FREQUENT

1. `a1 = aggregate(s=[], v=token)`
2. `logits1 = project(inp=a1, op=(inp==out),
special_op=(uniform selection))`
3. `prediction = softmax(logits1)`



We present more decompiled programs in the paper

Task	Program
Binary Majority	Sec. J.1
Binary Majority Interleave	Sec. J.2 J.3
Count	Sec. J.4 J.5
Most Frequent	Sec. J.6 J.7
Sort	Sec. J.8
Unique Bigram Copy	Sec. J.9 J.10
Unique Copy	Sec. J.11
Unique Reverse	Sec. J.12 J.13
<hr/>	
D_{12}	Sec. K.1
D_2	Sec. K.2
D_4	Sec. K.3
$(aa)^*$	Sec. K.4
$aa^*bb^*cc^*dd^*ee^*$	Sec. K.5
$\{a, b\}^*d\{b, c\}^*$	Sec. K.6
Tomita1	Sec. K.7
Tomita2	Sec. K.8
Tomita7	Sec. K.9

Applicable only:

- for length-generalizable models that actually internally implement RASP algorithms

Applicable only:

- for length-generalizable models that actually internally implement RASP algorithms
- when layer norm can be linearized, however, empirically this is true for length-generalizing models

What about LLMs?

- LLMs as a whole likely cannot be decompiled into a short RASP program

What about LLMs?

- LLMs as a whole likely cannot be decompiled into a short RASP program
- However, they may implement short RASP subprograms

What about LLMs?

- LLMs as a whole likely cannot be decompiled into a short RASP program
- However, they may implement short RASP subprograms
- Future work!

Conclusion

Method for extracting simple (D-)RASP programs from trained transformers.

Conclusion

Method for extracting simple (D-)RASP programs from trained transformers.

Show that transformers trained on tasks expressible in (C-)RASP often implement interpretable (D-)RASP programs.

Thank you!