# Inverse Tracr

William Baker

# What is RASP

A programming language with expressive power equal to a transformer

Map - Applies a function to a vector

Sequence Map - Applies a function to 2 vectors

Select - applies a predicate operation to vectors K and Q over $k_i$ and $q_j$, forming matrix $S_{ij}$
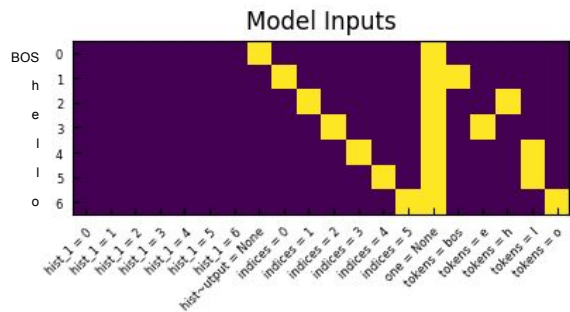
**Aggregate** - Given a boolean matrix and vector, use each row of the matrix to mask the vector and take the average over the non-masked values forming a column vector

```
e.g. aggregate(sel, [124]):
  F T T -> 1 * 0 + 2 * 1 + 4 * 1 / 2 = 3
  F F F -> 0 + 0 + 0 = 0
  T F F -> 1 * 1 + 0 + 0 / 1 = 1
  => [301]
```

**Select Width -** computes the average over rows of a matrix returning a vector
~ np.sum(axis=1)

# Tracr Overview

```
same_tok = rasp.Select(rasp.tokens, rasp.tokens,
                       rasp.Comparison.EQ).named("same_tok")
return rasp.SelectorWidth(same_tok).named("hist")
```
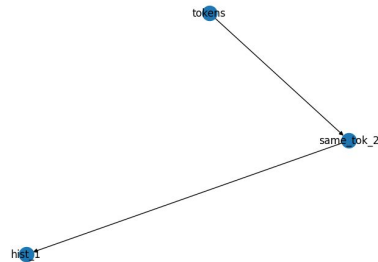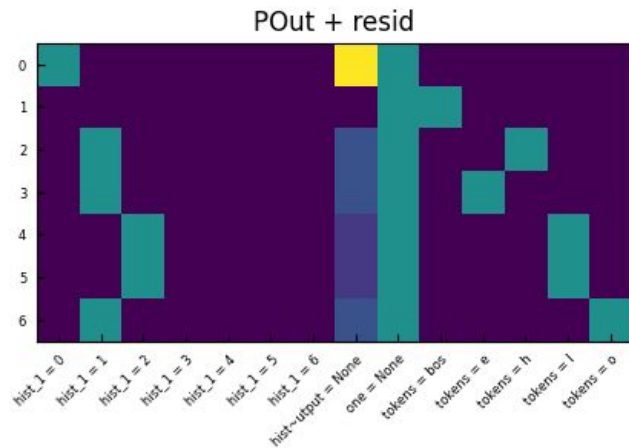
Computational Graph

Extract Craft Components

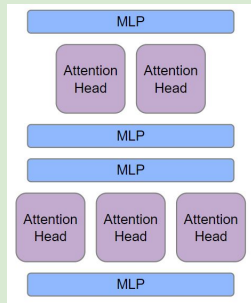Compile Craft Components into model

JAX model

Input = hello:

# Model

| Timestep | Terminal Flag | Param Block | Architecture |
|----------|:-------------:|:-----------:|:------------:|
| fst | 1 | <PARAMS> | <ARCH> |
| snd | 1 | <PARAMS> | <ARCH> |
| w_qk | 0 | <PARAMS> | <ARCH> |
| w_qk | 1 | <PARAMS> | <ARCH> |
| w_ov | 1 | <PARAMS> | <ARCH> |
| <PAD> | | | |
| PROG_START | <PAD> | | |
| <PAD> | | | |
| PROG_END | <PAD> | | |

Embed → GPT2 LARGE/ MEDIUM → Un-Embed

Ignored for Causal Masking

| Operator | Arg 1 | Arg 2 | Arg 3 | Return Var |
|----------|:-----:|:-----:|:-----:|:----------:|
| PROG_START | <PAD> | | | |
| MAP | LAM_ADD | indices | NA | v1 |
| Sequence Map | LAM_EQ | tokens | v1 | v2 |
| Sel | LAM_EQ | v1 | v2 | v3 |
| Aggregate | v3 | v1 | NA | v4 |
| PROG_END | <PAD> | | | |
| <PAD> | | | | |

# Initial Results

125M          377M          Tripled dataset size          774M          Added more variability in Residual stream size          774M

# Next steps

Try distilled models that will generalise beyond tracr

# Transformer Parameter Invariance

```python
def apply(self, x: bases.VectorInBasis) -> bases.VectorInBasis:
  assert x in self.residual_space
  # seq_len x query_space
  queries = x.project(self.w_qk.left_space)
  # seq_len x key_space
  keys = x.project(self.w_qk.right_space)

  attn_matrix = queries.magnitudes @ self.w_qk.matrix @ keys.magnitudes.T

  if self.causal:
    # The 1 gives us the matrix above the diagonal.
    mask = np.triu(np.full_like(attn_matrix, -np.inf), 1)
    attn_matrix = attn_matrix + mask

  attn_weights = _np_softmax(attn_matrix)  # seq_len_from, seq_len_to
  values = self.w_ov_residual(x).magnitudes  # seq_len_to, d_model

  magnitudes = attn_weights @ values  # seq_len_from, d_model
  return bases.VectorInBasis(sorted(self.residual_space.basis), magnitudes)
```

# Next steps

Test meta model on a subset of BERT parameters?