# Project 4: Pretraining

SW10 – Jannine Meier – FS24

# Preprocessing: Winogrande dataset

"I moved the couch from the garage to the backyard to create space. The ⎵ is small."

**1. Replace the blank with each option and combine the two sequences**

"... to create space. The couch is small."   + [SEP] +   "... to create space. The garage is small."

🚫 *Stemming and Lemmatizing*    **2. Tokenize the sequence using BertTokenizer**    *Remove punctuation and stopwords* 🚫
-> lowercase, add special tokens, add padding to max_length

[CLS] i moved ... to create space . the couch is small . [SEP] i
moved ... to create space . the garage is small . [SEP] [PAD] ... [PAD]

**3. Return for each sequence an input_id, attention_mask and a label**
label remapping to right option

- **Input_ids**: unique ID for each token
- **Attention_mask**: 1 for tokens, 0 for padding
- **Label**: 0 if first sentence is true, 1 otherwise

# Preprocessing: Anagram datasets

**1. Replace the <sep> with the Bert seperator token**

b p k <sep> k p b

b p k [SEP] k p b

**2. Tokenize the sequence using BertTokenizer**
-> add special tokens, add padding to max_length

[CLS] b p k [SEP] k p b [SEP] [PAD] ... [PAD]

**3. Return for each sequence an input_id, attention_mask and a label**

- **Input_ids:** unique ID for each token
- **Attention_mask:** 1 for tokens, 0 for padding
- **Label:** 1 if it is an anagram, 0 otherwise

# Model: Pretrained encoder transformer

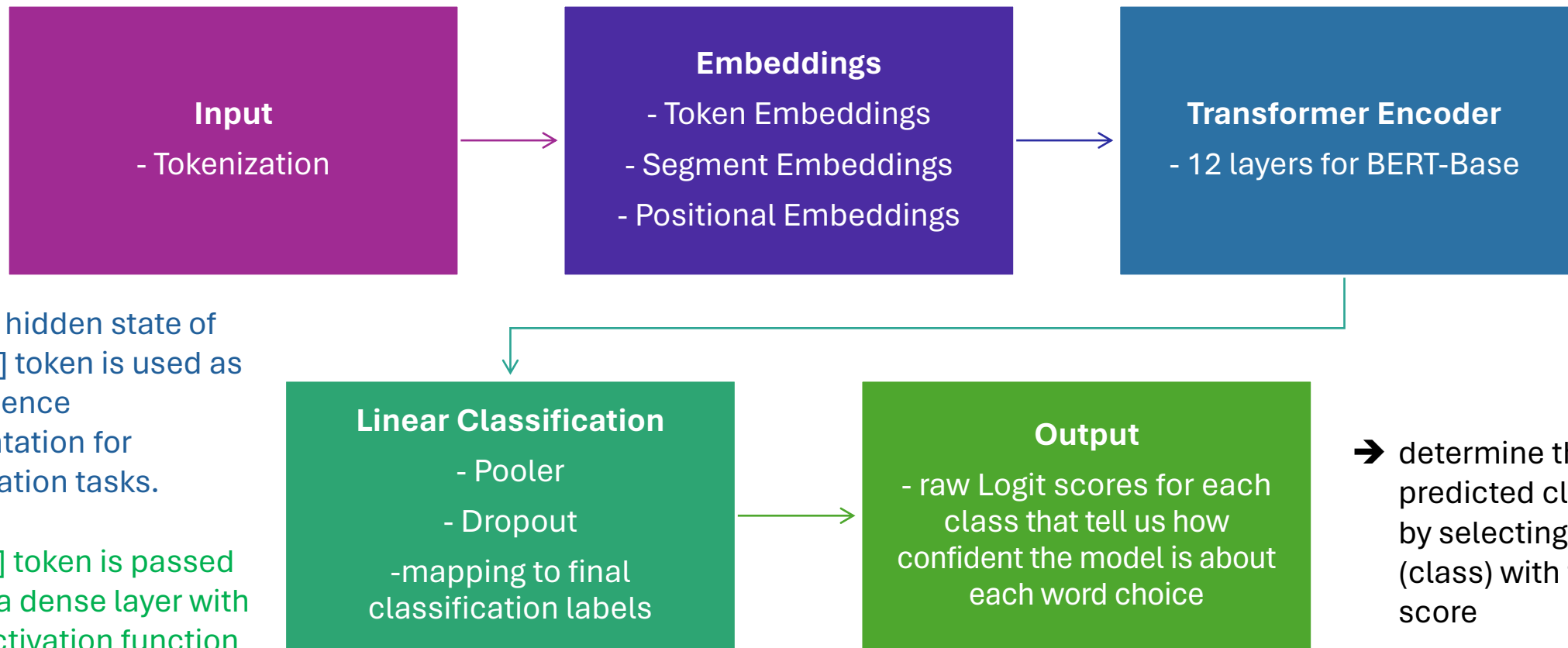## BertForSequenceClassification

(a pre-trained BERT model with a classification layer on top)

**Model Specification**

- Pre-trained Knowledge
    - great language understanding

- Bidirectional Contextual Understanding
    - consider both the words before and after the target word in the sentence

- Fine-Tuning for Our Specific Need
    - classify sequences into categories
    - decide if a sentence makes sense with one word or another

# Model: Pretrained encoder transformer

## BertForSequenceClassification

**Input**

- Tokenization

**Embeddings**

- Token Embeddings
- Segment Embeddings
- Positional Embeddings

**Transformer Encoder**

- 12 layers for BERT-Base

The final hidden state of the [CLS] token is used as the sequence representation for classification tasks.

The [CLS] token is passed through a dense layer with a tanh activation function to create a pooled output.

**Linear Classification**

- Pooler
- Dropout
- mapping to final classification labels

**Output**

- raw Logit scores for each class that tell us how confident the model is about each word choice

➔ determine the predicted class label by selecting the index (class) with the highest score

# Experiment settings

**Sweep configurations**

- method: grid
- learning_rate: 1e-6, 1e-5, 1e-4, 1e-3, 1e-2

**Fixed parameters**

- batch_size: 32/64 (max possible)
- num_epochs: 100/10
- warmup_ratio: 0.1
- early_stopping_patience: 20
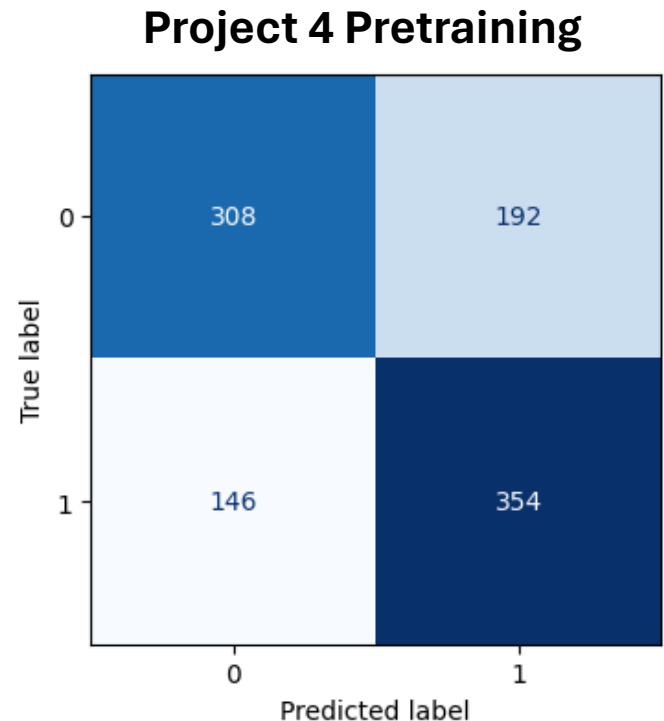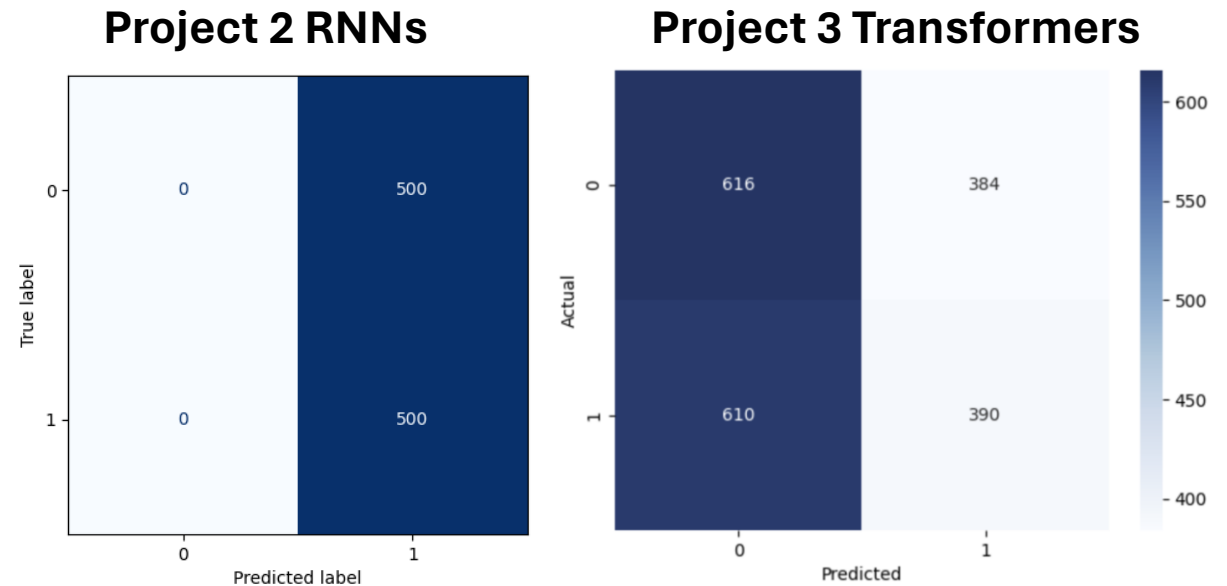- Loss Function: CrossEntropy
- Optimizer: AdamW with weight_decay: 0.1

# Results: Winogrande

**Best model**

- Validation accuracy: 55.3%
  - last project 50.9%
- Test accuracy: 66.2%
  - last project 50.0%
- Small increase in accuracy
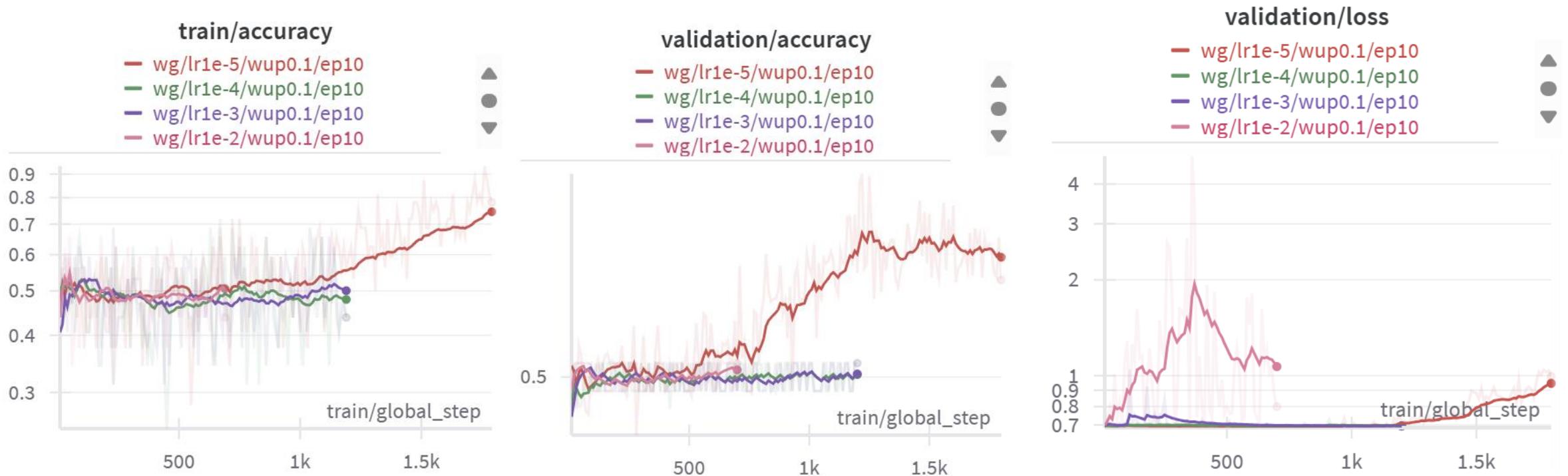- Slightly better performance than random

**Confusion-matrix**

- more balance in predictions



**Project 2 RNNs**

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 0 | 500 |
| True 1 | 0 | 500 |

**Project 3 Transformers**

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | 616 | 384 |
| Actual 1 | 610 | 390 |

**Project 4 Pretraining**

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| True 0 | 308 | 192 |
| True 1 | 146 | 354 |

# Results: Winogrande

**Finally some learning progress in train and validation!**
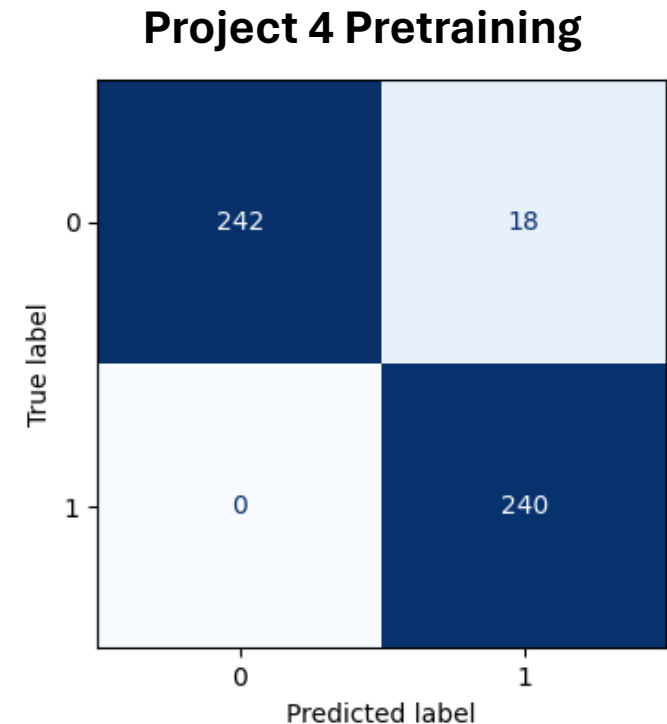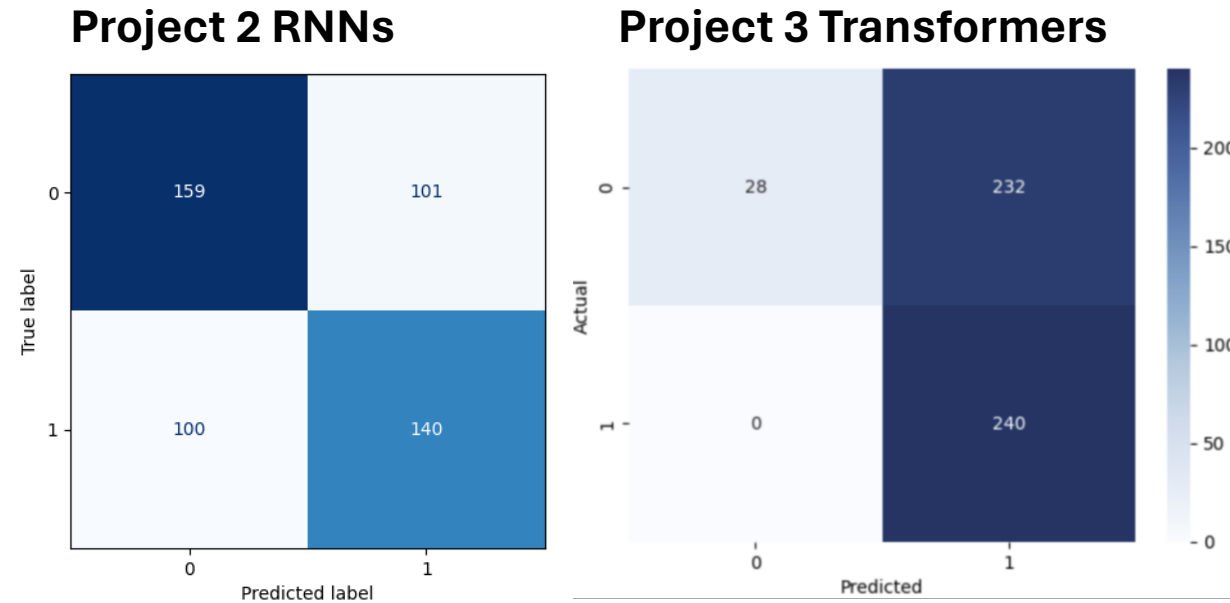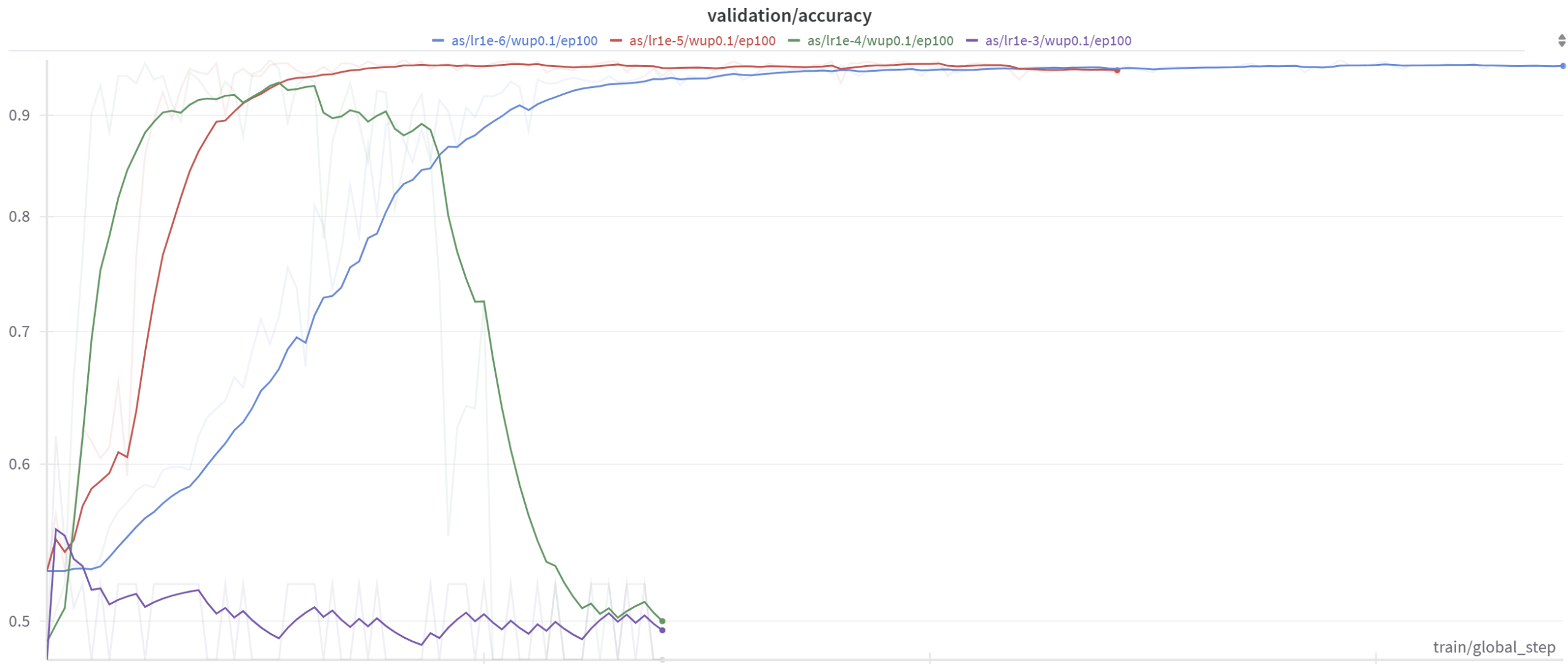
# Results:
# Anagram Small

## Best model

- Validation accuracy: 95.6%
  - last project 56.4%

- Test accuracy: 96.4%
  - last project 53.6%

## Confusion-matrix

- more balance in predictions
- Performs much better than previous projects



**Project 2 RNNs**

**Project 3 Transformers**

**Project 4 Pretraining**

# Results: Anagram Small



validation/accuracy

as/lr1e-6/wup0.1/ep100 — as/lr1e-5/wup0.1/ep100 — as/lr1e-4/wup0.1/ep100 — as/lr1e-3/wup0.1/ep100
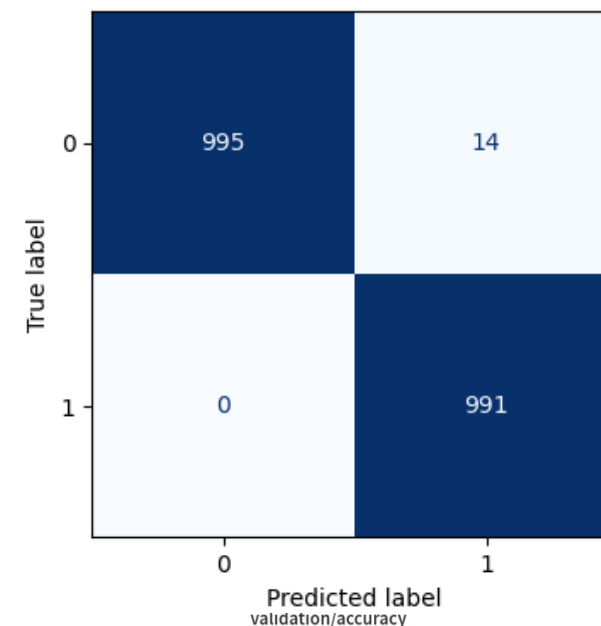
train/global_step

# Results:
# Anagram Large

**Best model**

- Validation accuracy: 99.5%
  - last project 99.9%

- Test accuracy: 99.3%
  - last project 97.2%

**Confusion-matrix**

- Almost perfect performance

# Conclusions

**Interpretation**

**Winogrande dataset**

- still too complex for my model choice

**Anagram datasets**

- almost perfect performance

➔ For all tasks small learning rates (1e-5 and 1e-6) showed best results

**Lessons Learned**

**Less Data required**

- with the pretrained transformer I already got good results on the smaller dataset as well

**Warmup learning rate**

- try fixed steps and ratio