

# Libraries and Tools

# Transformers, AllenNLP

LING575 Analyzing Neural Language Models

Shane Steinert-Threlkeld

Apr 27 2022

# Outline

- Very helpful tools
  - 🤗 Transformers
  - AllenNLP
    - Walk-through of a classifier and a tagger
- Second half: tips/tricks for experiment running and paper writing

# Transformers

<https://huggingface.co/transformers>

# Where to get LMs to analyze?

- RNNs: see week 3 slides
  - Josefewicz et al “Exploring the limits...”
  - Gulordava et al “Colorless green ideas...”
  - ELMo via AllenNLP (about which more later)
- Effectively a unique API for each model
- All (essentially) Transformer-based models: HuggingFace!

# Overview of the Library

- Access to many variants of many very large LMs (BERT, RoBERTa, XLNET, ALBERT, T5, language-specific models, ...) with fairly consistent API
  - Build tokenizer + model from string for name or config
  - Then use just like any PyTorch nn.Module
- Emphasis on ease-of-use
  - E.g. low barrier-to-entry to *using* the models, including for analysis
  - [new `pipeline` abstraction too, but I think this is *too easy* for most analysis / probing purposes, but can work if all you need are model judgments on data]
- Interoperable with PyTorch or TensorFlow 2.0

# Example: Tokenization

```
>>> from transformers import AutoTokenizer
```

```
>>> tokenizer = AutoTokenizer.from_pretrained("bert-base-cased")
```

```
>>> encoded_input = tokenizer("Do not meddle in the affairs of wizards, for they are subtle and quick to  
>>> print(encoded_input)  
{'input_ids': [101, 2079, 2025, 19960, 10362, 1999, 1996, 3821, 1997, 16657, 1010, 2005, 2027, 2024, 1125],  
  'token_type_ids': [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
  'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]}
```

# Example: Tokenization

```
>>> from transformers import AutoTokenizer
```

```
>>> tokenizer = AutoTokenizer.from_pretrained("bert-base-cased")
```

```
>>> encoded_input = tokenizer("Do not meddle in the affairs of wizards, for they are subtle and quick to  
>>> print(encoded_input)  
{'input_ids': [101, 2079, 2025, 19960, 10362, 1999, 1996, 3821, 1997, 16657, 1010, 2005, 2027, 2024, 1125],  
  'token_type_ids': [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
  'attention_mask': [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]}
```

```
>>> tokenizer.decode(encoded_input["input_ids"])  
'[CLS] Do not meddle in the affairs of wizards, for they are subtle and quick to anger. [SEP]'
```

# Example: Tokenizing a Batch

```
>>> batch_sentences = [  
...     "But what about second breakfast?",  
...     "Don't think he knows about second breakfast, Pip.",  
...     "What about elevensies?",  
... ]  
>>> encoded_input = tokenizer(batch_sentences, padding=True)  
>>> print(encoded_input)  
{  
  'input_ids': [[101, 1252, 1184, 1164, 1248, 6462, 136, 102, 0, 0, 0, 0, 0, 0, 0],  
                [101, 1790, 112, 189, 1341, 1119, 3520, 1164, 1248, 6462, 117, 21902, 1643, 119, 102],  
                [101, 1327, 1164, 5450, 23434, 136, 102, 0, 0, 0, 0, 0, 0, 0, 0]],  
  'token_type_ids': [[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
                     [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
                     [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]],  
  'attention_mask': [[1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0],  
                    [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],  
                    [1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0]]  
}
```



# Example: Tokenizing a Batch

```
>>> batch_sentences = [  
...     "But what about second breakfast?",  
...     "Don't think he knows about second breakfast, Pip.",  
...     "What about elevensies?",  
... ]  
>>> encoded_input = tokenizer(batch_sentences, padding=True)  
>>> print(encoded_input)  
{'input_ids': [[101, 1252, 1184, 1164, 1248, 6462, 136, 102, 0, 0, 0, 0, 0, 0, 0],  
               [101, 1790, 112, 189, 1341, 1119, 3520, 1164, 1248, 6462, 117, 21902, 1643, 119, 102],  
               [101, 1327, 1164, 5450, 23434, 136, 102, 0, 0, 0, 0, 0, 0, 0, 0]],  
 'token_type_ids': [[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
                    [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],  
                    [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]],  
 'attention_mask': [[1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0],  
                   [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1],  
                   [1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0]]}
```

Add ``return_tensors="pt"`` to get these outputs as PyTorch Tensors

# Example: Forward Pass

```
>>> from transformers import BertTokenizer, BertModel
>>> import torch

>>> tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
>>> model = BertModel.from_pretrained("bert-base-uncased")

>>> inputs = tokenizer("Hello, my dog is cute", return_tensors="pt")
>>> outputs = model(**inputs)

>>> last_hidden_states = outputs.last_hidden_state
```

# Outputs from the forward pass

- Outputs are *usually* Python objects with various attributes corresponding to different model outputs (NB: can be *tuples of Tensors* if specified, but I recommend against that)
- BERT, by default, gives two things:
  - **last\_hidden\_state**: sequence of hidden states at the last layer of the model.  
Shape: (batch\_size, max\_length, embedding\_dimension)
  - **pooler\_output**: embedding of '[CLS]' token, passed through one tanh layer (more on this later)  
Shape: (batch\_size, embedding\_dimension)

# Getting more out of a model

```
from transformers import BertModel
model = BertModel.from_pretrained("bert-base-uncased")

outputs = model(inputs, output_hidden_states=True,
output_attentions=True)
```

# Getting more out of a model

```
from transformers import BertModel
model = BertModel.from_pretrained("bert-base-uncased")

outputs = model(inputs, output_hidden_states=True,
output_attentions=True)
```

- Now, the output object has additional attributes:
  - **hidden\_states**: A tuple of tensors, one for each layer. Length: # layers  
Shape of each: (batch\_size, max\_length, embedding\_dimension)
  - **attentions**: tuple of tensors, one for each layer. Length: # layers  
Shape of each: (batch\_size, num\_heads, max\_length, max\_length)
- [Can also be done with BertConfig object]

# What the library does well

- Very easy tokenization
- Forward pass of models
  - Exposing as many internals as possible
    - All layers, attention heads, etc
- As unified an interface as possible
  - But: different models have different properties, controlled by Configs or by arguments
  - Read the docs carefully!
    - e.g. [https://huggingface.co/docs/transformers/model\\_doc/bert#transformers.BertModel](https://huggingface.co/docs/transformers/model_doc/bert#transformers.BertModel)
    - The docs for `forward` just below explain the inputs/outputs



# More Info

- Model search: <https://huggingface.co/models>
- Dataset search: <https://huggingface.co/datasets>
- Relatively new portions of the library (Trainer) may be useful for probing, but we'll look at another route for that now.

# AllenNLP

<https://allenai.org/allennlp/software/allennlp-library>



# Overview of AllenNLP

- Built on top of PyTorch
- Flexible data API
- Abstractions for common use cases in NLP
  - e.g. take a sequence of representations and give me a single one
- Modular:
  - Because of that, can swap in and out different options, for good experiments
- Declarative model-building / training via config files
- See <https://github.com/allenai/writing-code-for-nlp-research-emnlp2018>
- Guide: <https://guide.allennlp.org/> <— very helpful for explaining some of the abstractions

# Some Advantages

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:
- Training loop:

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:
- Training loop: 

```
for each epoch:  
  for each batch:  
    get model outputs on batch  
    compute loss  
    compute gradients  
    update parameters
```

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:
- Training loop: 

```
for each epoch:  
    for each batch:  
        get model outputs on batch  
        compute loss  
        compute gradients  
        update parameters
```
- Not *that* complicated, but:

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:
- Training loop: 

```
for each epoch:  
    for each batch:  
        get model outputs on batch  
        compute loss  
        compute gradients  
        update parameters
```
- Not *that* complicated, but:
  - Early stopping

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:
- Training loop: 

```
for each epoch:  
    for each batch:  
        get model outputs on batch  
        compute loss  
        compute gradients  
        update parameters
```
- Not *that* complicated, but:
  - Early stopping
  - Check-pointing (saving best model(s))



# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:

- Training loop: 

```
for each epoch:  
    for each batch:  
        get model outputs on batch  
        compute loss  
        compute gradients  
        update parameters
```

- Not *that* complicated, but:
  - Early stopping
  - Check-pointing (saving best model(s))
  - Generating and padding the batches

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:

- Training loop: 

```
for each epoch:
    for each batch:
        get model outputs on batch
        compute loss
        compute gradients
        update parameters
```

- Not *that* complicated, but:
  - Early stopping
  - Check-pointing (saving best model(s))
  - Generating and padding the batches
  - Logging results

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:

- Training loop: 

```
for each epoch:
    for each batch:
        get model outputs on batch
        compute loss
        compute gradients
        update parameters
```

- Not *that* complicated, but:
  - Early stopping
  - Check-pointing (saving best model(s))
  - Generating and padding the batches
  - Logging results
  - ....

# Some Advantages

- Focus on modeling / experimenting, not writing boilerplate, e.g.:

- Training loop: 

```
for each epoch:
    for each batch:
        get model outputs on batch
        compute loss
        compute gradients
        update parameters
```

- Not *that* complicated, but:

- Early stopping
- Check-pointing (saving best model(s))
- Generating and padding the batches
- Logging results

- .... `allennlp train myexperiment.jsonnet`

# Example Abstractions

- TextFieldEmbedder
  - Seq2SeqEncoder
  - Seq2VecEncoder
  - Attention
  - ...
- 
- Allows for easy swapping of different choices at every level in your model.

# AllenNLP Bert Example

- See <https://github.com/shanest/allennlp-bert-example> [linked on course webpage as well]
- Using AllenNLP to probe BERT for two tasks:
  - Classification [[Stanford Sentiment Treebank](#)]
  - Tagging [[Semantic Tagging](#)]

# Classifying

# Overall Structure (Classification)

```
model_string = "bert-base-uncased"

tokenizer = PretrainedTransformerTokenizer(model_string)
token_indexer = PretrainedTransformerIndexer(model_string)

reader = SSTDatasetReader(tokenizer, {"tokens": token_indexer})
train_path = "sst/trees/train.txt"
dev_path = "sst/trees/dev.txt"

train_dataset = reader.read(train_path)
val_dataset = reader.read(dev_path)

print(list(train_dataset)[0])

vocab = Vocabulary.from_instances(chain(train_dataset, val_dataset))

bert_token_embedder = PretrainedTransformerEmbedder(model_string)
bert_textfield_embedder = BasicTextFieldEmbedder({"tokens": bert_token_embedder})
cls_pooler = ClsPooler(bert_token_embedder.get_output_dim())

model = BertClassifier(
    vocab, bert_textfield_embedder, cls_pooler, freeze_encoder=False
)

data_loader = MultiProcessDataLoader(reader, train_path, batch_size=32)
data_loader.index_with(vocab)

trainer = GradientDescentTrainer(
    model=model,
    optimizer=optim.Adam(model.parameters()),
    serialization_dir="/tmp/test",
    data_loader=data_loader,
    train_dataset=train_dataset,
    validation_dataset=val_dataset,
    patience=5,
    num_epochs=30,
)
trainer.train()
```



# Overall Structure (Classification)

```
model_string = "bert-base-uncased"

tokenizer = PretrainedTransformerTokenizer(model_string)
token_indexer = PretrainedTransformerIndexer(model_string)

reader = SSTDatasetReader(tokenizer, {"tokens": token_indexer})
train_path = "sst/trees/train.txt"
dev_path = "sst/trees/dev.txt"

train_dataset = reader.read(train_path)
val_dataset = reader.read(dev_path)

print(list(train_dataset)[0])

vocab = Vocabulary.from_instances(chain(train_dataset, val_dataset))

bert_token_embedder = PretrainedTransformerEmbedder(model_string)
bert_textfield_embedder = BasicTextFieldEmbedder({"tokens": bert_token_embedder})
cls_pooler = ClsPooler(bert_token_embedder.get_output_dim())

model = BertClassifier(
    vocab, bert_textfield_embedder, cls_pooler, freeze_encoder=False
)

data_loader = MultiProcessDataLoader(reader, train_path, batch_size=32)
data_loader.index_with(vocab)

trainer = GradientDescentTrainer(
    model=model,
    optimizer=optim.Adam(model.parameters()),
    serialization_dir="/tmp/test",
    data_loader=data_loader,
    train_dataset=train_dataset,
    validation_dataset=val_dataset,
    patience=5,
    num_epochs=30,
)
trainer.train()
```

← DatasetReader

# Overall Structure (Classification)

```
model_string = "bert-base-uncased"

tokenizer = PretrainedTransformerTokenizer(model_string)
token_indexer = PretrainedTransformerIndexer(model_string)

reader = SSTDatasetReader(tokenizer, {"tokens": token_indexer})
train_path = "sst/trees/train.txt"
dev_path = "sst/trees/dev.txt"

train_dataset = reader.read(train_path)
val_dataset = reader.read(dev_path)

print(list(train_dataset)[0])

vocab = Vocabulary.from_instances(chain(train_dataset, val_dataset))

bert_token_embedder = PretrainedTransformerEmbedder(model_string)
bert_textfield_embedder = BasicTextFieldEmbedder({"tokens": bert_token_embedder})
cls_pooler = ClsPooler(bert_token_embedder.get_output_dim())

model = BertClassifier(
    vocab, bert_textfield_embedder, cls_pooler, freeze_encoder=False
)

data_loader = MultiProcessDataLoader(reader, train_path, batch_size=32)
data_loader.index_with(vocab)

trainer = GradientDescentTrainer(
    model=model,
    optimizer=optim.Adam(model.parameters()),
    serialization_dir="/tmp/test",
    data_loader=data_loader,
    train_dataset=train_dataset,
    validation_dataset=val_dataset,
    patience=5,
    num_epochs=30,
)
trainer.train()
```

← DatasetReader

← Model

# Overall Structure (Classification)

```
model_string = "bert-base-uncased"

tokenizer = PretrainedTransformerTokenizer(model_string)
token_indexer = PretrainedTransformerIndexer(model_string)

reader = SSTDatasetReader(tokenizer, {"tokens": token_indexer})
train_path = "sst/trees/train.txt"
dev_path = "sst/trees/dev.txt"

train_dataset = reader.read(train_path)
val_dataset = reader.read(dev_path)

print(list(train_dataset)[0])

vocab = Vocabulary.from_instances(chain(train_dataset, val_dataset))

bert_token_embedder = PretrainedTransformerEmbedder(model_string)
bert_textfield_embedder = BasicTextFieldEmbedder({"tokens": bert_token_embedder})
cls_pooler = ClsPooler(bert_token_embedder.get_output_dim())

model = BertClassifier(
    vocab, bert_textfield_embedder, cls_pooler, freeze_encoder=False
)

data_loader = MultiProcessDataLoader(reader, train_path, batch_size=32)
data_loader.index_with(vocab)

trainer = GradientDescentTrainer(
    model=model,
    optimizer=optim.Adam(model.parameters()),
    serialization_dir="/tmp/test",
    data_loader=data_loader,
    train_dataset=train_dataset,
    validation_dataset=val_dataset,
    patience=5,
    num_epochs=30,
)
trainer.train()
```

← DatasetReader

← Model

← Data Loader/Iterator

# Overall Structure (Classification)

```
model_string = "bert-base-uncased"

tokenizer = PretrainedTransformerTokenizer(model_string)
token_indexer = PretrainedTransformerIndexer(model_string)

reader = SSTDatasetReader(tokenizer, {"tokens": token_indexer})
train_path = "sst/trees/train.txt"
dev_path = "sst/trees/dev.txt"

train_dataset = reader.read(train_path)
val_dataset = reader.read(dev_path)

print(list(train_dataset)[0])

vocab = Vocabulary.from_instances(chain(train_dataset, val_dataset))

bert_token_embedder = PretrainedTransformerEmbedder(model_string)
bert_textfield_embedder = BasicTextFieldEmbedder({"tokens": bert_token_embedder})
cls_pooler = ClsPooler(bert_token_embedder.get_output_dim())

model = BertClassifier(
    vocab, bert_textfield_embedder, cls_pooler, freeze_encoder=False
)

data_loader = MultiProcessDataLoader(reader, train_path, batch_size=32)
data_loader.index_with(vocab)

trainer = GradientDescentTrainer(
    model=model,
    optimizer=optim.Adam(model.parameters()),
    serialization_dir="/tmp/test",
    data_loader=data_loader,
    train_dataset=train_dataset,
    validation_dataset=val_dataset,
    patience=5,
    num_epochs=30,
)
trainer.train()
```

← DatasetReader

← Model

← Data Loader/Iterator

← Trainer

# Basic Components: Dataset Reader

- Datasets are collections of *Instances*, which are collections of *Fields*
  - For text classification, e.g.: one TextField, one LabelField
  - Many more: <https://guide.allennlp.org/reading-data>
- DatasetReaders..... read data sets. Two primary methods:
  - `_read(file)`: reads data from disk, yields Instances. By calling:
    - `text_to_instance` (variable signature)
      - Processing of the “raw” data from disk into final form
      - Produces one Instance at a time



# DatasetReader: Stanford Sentiment Treebank

- One line from train.txt:

(3 (2 (2 The) (2 Rock)) (4 (3 (2 is) (4 (2 destined) (2 (2 (2 (2 to) (2 (2 be) (2 (2 the) (2 (2 21st) (2 (2 (2 Century) (2 's)) (2 (3 new) (2 (2 ``) (2 Conan)))))))) (2 ") (2 and)) (3 (2 that) (3 (2 he) (3 (2 's) (3 (2 going) (3 (2 to) (4 (3 (2 make) (3 (3 (2 a) (3 splash)) (2 (2 even) (3 greater)))) (2 (2 than) (2 (2 (2 (1 (2 Arnold) (2 Schwarzenegger)) (2 ,)) (2 (2 Jean-Claud) (2 (2 Van) (2 Damme)))) (2 or)) (2 (2 Steven) (2 Segal)))))))))) (2 .)))

- Core of \_read:

```
parsed_line = Tree.fromstring(line)
instance = self.text_to_instance(parsed_line.leaves(), parsed_line.label())
if instance is not None:
    yield instance
```

- Core of text\_to\_instance:

```
if self._tokenizer:
    new_tokens = self._tokenizer.tokenize(' '.join(tokens))
else:
    new_tokens = [Token(token) for token in tokens]
text_field = TextField(new_tokens, token_indexers=self._token_indexers)
fields: Dict[str, Field] = {"tokens": text_field}
```

...

```
fields["label"] = LabelField(sentiment)
return Instance(fields)
```

# Model

```
@Model.register("bert_classifier")
class BertClassifier(Model):
    def __init__(
        self,
        vocab: Vocabulary,
        embedder: TextFieldEmbedder,
        pooler: Seq2VecEncoder,
        freeze_encoder: bool = True,
    ) -> None:
        super().__init__(vocab)

        self.vocab = vocab
        self.embedder = embedder
        self.pooler = pooler
        self.freeze_encoder = freeze_encoder

        for parameter in self.embedder.parameters():
            parameter.requires_grad = not self.freeze_encoder

        in_features = self.pooler.get_output_dim()
        out_features = vocab.get_vocab_size(namespace="labels")

        self._classification_layer = torch.nn.Linear(in_features, out_features)
        self._accuracy = CategoricalAccuracy()
        self._loss = torch.nn.CrossEntropyLoss()
```

# Model

```
@Model.register("bert_classifier")
class BertClassifier(Model):
    def __init__(
        self,
        vocab: Vocabulary,
        embedder: TextFieldEmbedder,
        pooler: Seq2VecEncoder,
        freeze_encoder: bool = True,
    ) -> None:
        super().__init__(vocab)

        self.vocab = vocab
        self.embedder = embedder
        self.pooler = pooler
        self.freeze_encoder = freeze_encoder

        for parameter in self.embedder.parameters():
            parameter.requires_grad = not self.freeze_encoder

        in_features = self.pooler.get_output_dim()
        out_features = vocab.get_vocab_size(namespace="labels")

        self._classification_layer = torch.nn.Linear(in_features, out_features)
        self._accuracy = CategoricalAccuracy()
        self._loss = torch.nn.CrossEntropyLoss()
```

Fine tune or not





# Model

```
def forward( # type: ignore
    self, tokens: Dict[str, torch.Tensor], label: torch.IntTensor = None
) -> Dict[str, torch.Tensor]:

    # (batch_size, max_len, embedding_dim)
    embeddings = self.embedder(tokens)

    # get the pooled representation of the tokens in each sentence
    # e.g. [CLS] rep, mean pool, ...
    # (batch_size, embedding_dim)
    sentence_representation = self.pooler(embeddings)

    # apply classification layer
    # (batch_size, num_labels)
    logits = self._classification_layer(sentence_representation)

    probs = torch.nn.functional.softmax(logits, dim=-1)

    output_dict = {"logits": logits, "probs": probs}

    if label is not None:
        loss = self._loss(logits, label.long().view(-1))
        output_dict["loss"] = loss
        self._accuracy(logits, label)

    return output_dict
```

← NB: frozen embeddings can be pre-computed for efficiency

# Where was BERT?

- In the TextFieldEmbedder!
- In run\_classifying.py: initialized a PretrainedTransformerEmbedder
  - AllenNLP has wrappers around HuggingFace
  - But note: to extract more from a model, you'll probably need to write your own class, using the existing ones as inspiration

# Config file (classifying\_experiment.jsonnet)

```
local bert_model = "bert-base-uncased";
{
  "dataset_reader": {
    "type": "sst_reader",
    "tokenizer": {
      "type": "pretrained_transformer",
      "model_name": bert_model,
    },
    "token_indexers": {
      "tokens": {
        "type": "pretrained_transformer",
        "model_name": bert_model,
      }
    }
  },
  "train_data_path": "sst/trees/train.txt",
  "validation_data_path": "sst/trees/dev.txt",
}
```

@DatasetReader.register("sst\_reader")

Arguments to SSTReader!

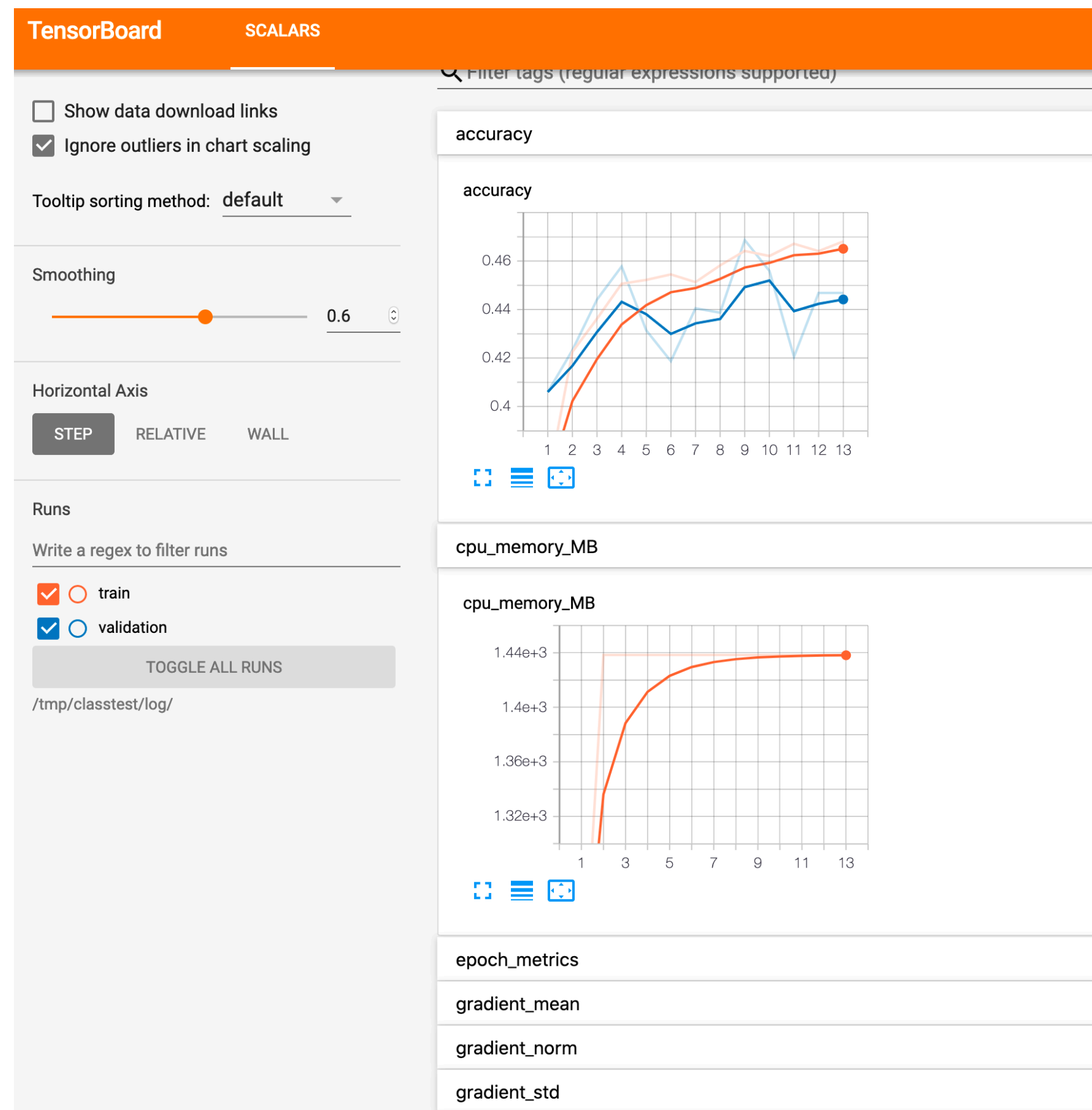
# Config file (classifying\_experiment.jsonnet)

```
"model": {
  "type": "bert_classifier",
  "embedder": {
    "type": "basic",
    "token_embedders": {
      "tokens": {
        "type": "pretrained_transformer",
        "model_name": bert_model
      }
    }
  },
  "pooler": {
    # NB: for probing, cls_pooler and boe_pooler are good choices
    # bert_pooler actually does more than what is wanted in that scenario
    "type": "cls_pooler",
    "embedding_dim": 768,
  },
  "freeze_encoder": true,
},
"data_loader": {
  "batch_size": 32
},
"trainer": {
  "optimizer": {
    "type": "adam",
    "lr": 0.001
  },
  "validation_metric": "+accuracy",
  "checkpointer": {
    "keep_most_recent_by_count": 1
  },
  "num_epochs": 30,
  "grad_norm": 10.0,
  "patience": 5,
  "cuda_device": -1
}
```

```
allennlp train classifying_experiment.jsonnet \
  --serialization-dir test \
  --include-package classifying
```

# TensorBoard

```
tensorboard --logdir /serialization_dir/log
```



Use SSH port forwarding to  
view server-side results locally

# Tagging

# Tagging



# Tagging

- The repository also has an example of training a semantic tagger
- Like POS tagging, but with a richer set of “semantic” tags

anaphoric	ANA	PRO	anaphoric & deictic pronouns: <i>he, she, I, him</i>	NOT	negation: <i>not, no, neither, without</i>	MOD	modality
		DEF	definite: <i>the, lo<sup>IT</sup>, der<sup>DE</sup></i>	NEC	necessity: <i>must, should, have to</i>		
		HAS	possessive pronoun: <i>my, her</i>	POS	possibility: <i>might, could, perhaps, alleged, can</i>		
		REF	reflexive & reciprocal pron.: <i>herself, each_other</i>				
		EMP	emphasizing pronouns: <i>himself</i>				
speech act	ACT	GRE	greeting & parting: <i>hi, bye</i>	SUB	subordinate relations: <i>that, while, because</i>	DSC	discourse
		ITJ	interjections, exclamations: <i>alas, ah</i>	COO	coordinate relations: <i>so, {, }, {; }, and</i>		
		HES	hesitation: <i>err</i>	APP	appositional relations: <i>{, }, which, {( }, {—}</i>		
		QUE	interrogative: <i>who, which, ?</i>	BUT	contrast: <i>but, yet</i>		
attribute	ATT	QUC	*concrete quantity: <i>two, six_million, twice</i>	PER	person: <i>Axl_Rose, Sherlock_Holmes</i>	NAM	named entity
		QUV	*vague quantity: <i>millions, many, enough</i>	GPE	geo-political entity: <i>Paris, Japan</i>		
		COL	*colour: <i>red, crimson, light_blue, chestnut_brown</i>	GPO	*geo-political origin: <i>Parisian, French</i>		
		IST	intersective: <i>open, vegetarian, quickly</i>	GEO	geographical location: <i>Alps, Nile</i>		
		SST	subsective: <i>skillful surgeon, tall kid</i>	ORG	organization: <i>IKEA, EU</i>		
		PRI	privative: <i>former, fake</i>	ART	artifact: <i>iOS_7</i>		
		DEG	*degree: <i>2 meters tall, 20 years old</i>	HAP	happening: <i>Eurovision_2017</i>		
		INT	intensifier: <i>very, much, too, rather</i>	UOM	unit of measurement: <i>meter, \$, %, degree_Celsius</i>		
		REL	relation: <i>in, on, 's, of, after</i>	CTC	*contact information: <i>112, info@mail.com</i>		
		SCO	score: <i>3-0, grade A</i>	URL	URL: <i>http://pmb.let.rug.nl</i>		
comparative	COM	EQU	equative: <i>as tall as John, whales are mammals</i>	LIT	*literal use of names: <i>his name is John</i>		
		MOR	comparative positive: <i>better, more</i>	NTH	*other names: <i>table 1a, equation (1)</i>		
		LES	comparative negative: <i>less, worse</i>	EXS	untensed simple: <i>to walk, is eaten, destruction</i>	EVE	events
		TOP	superlative positive: <i>most, mostly</i>	ENS	present simple: <i>we walk, he walks</i>		
		BOT	superlative negative: <i>worst, least</i>	EPS	past simple: <i>ate, went</i>		
		ORD	ordinal: <i>1st, 3rd, third</i>	EXG	untensed progressive: <i>is running</i>		
				EXT	untensed perfect: <i>has eaten</i>		
unnamed entity	UNE	CON	concept: <i>dog, person</i>	NOW	present tense: <i>is skiing, do ski, has skied, now</i>	TNS	tense & aspect
		ROL	role: <i>student, brother, prof., victim</i>	PST	past tense: <i>was baked, had gone, did go</i>		
		GRP	*group: <i>John {, } Mary and Sam gathered, a group of people</i>	FUT	future tense: <i>will, shall</i>		
deixis	DXS	DXP	*place deixis: <i>here, this, above</i>	PRG	*progressive: <i>has been being treated, aan_het<sup>NL</sup></i>		
		DXT	*temporal deixis: <i>just, later, tomorrow</i>	PFT	*perfect: <i>has been going/done</i>		
		DXD	*discourse deixis: <i>latter, former, above</i>				
logical	LOG	ALT	alternative & repetitions: <i>another, different, again</i>	DAT	*full date: <i>27.04.2017, 27/04/17</i>	TIM	temporal entity
		XCL	exclusive: <i>only, just</i>	DOM	day of month: <i>27th December</i>		
		NIL	empty semantics: <i>{.}, to, of</i>	YOC	year of century: <i>2017</i>		
		DIS	disjunction & exist. quantif.: <i>a, some, any, or</i>	DOW	day of week: <i>Thursday</i>		
		IMP	implication: <i>if, when, unless</i>	MOY	month of year: <i>April</i>		
		AND	conjunction & univ. quantif.: <i>every, and, who, any</i>	DEC	decade: <i>80s, 1990s</i>		
				CLO	clocktime: <i>8:45_pm, 10_o'clock, noon</i>		



# Tagging

- The repository also has an example of training a semantic tagger
- Like POS tagging, but with a richer set of “semantic” tags
- Issue: the data comes with its own tokenization:
  - BERT: ['the', 'ya', '##zuka', 'are', 'the', 'japanese', 'mafia', '.']

<b>ANA</b> anaphoric	<b>PRO</b> anaphoric & deictic pronouns: <i>he, she, I, him</i> <b>DEF</b> definite: <i>the, lo<sup>IT</sup>, der<sup>DE</sup></i> <b>HAS</b> possessive pronoun: <i>my, her</i> <b>REF</b> reflexive & reciprocal pron.: <i>herself, each_other</i> <b>EMP</b> emphasizing pronouns: <i>himself</i>	<b>NOT</b> negation: <i>not, no, neither, without</i> <b>NEC</b> necessity: <i>must, should, have to</i> <b>POS</b> possibility: <i>might, could, perhaps, alleged, can</i>	<b>MOD</b> modality
<b>ACT</b> speech act	<b>GRE</b> greeting & parting: <i>hi, bye</i> <b>ITJ</b> interjections, exclamations: <i>alas, ah</i> <b>HES</b> hesitation: <i>err</i> <b>QUE</b> interrogative: <i>who, which, ?</i>	<b>SUB</b> subordinate relations: <i>that, while, because</i> <b>COO</b> coordinate relations: <i>so, {, }, {; }, and</i> <b>APP</b> appositional relations: <i>{, }, which, {(), {—}</i> <b>BUT</b> contrast: <i>but, yet</i>	<b>DSC</b> discourse
<b>ATT</b> attribute	<b>QUC</b> *concrete quantity: <i>two, six_million, twice</i> <b>QUV</b> *vague quantity: <i>millions, many, enough</i> <b>COL</b> *colour: <i>red, crimson, light_blue, chestnut_brown</i> <b>IST</b> intersective: <i>open, vegetarian, quickly</i> <b>SST</b> subsective: <i>skillful surgeon, tall kid</i> <b>PRI</b> privative: <i>former, fake</i> <b>DEG</b> *degree: <i>2 meters tall, 20 years old</i> <b>INT</b> intensifier: <i>very, much, too, rather</i> <b>REL</b> relation: <i>in, on, 's, of, after</i> <b>SCO</b> score: <i>3-0, grade A</i>	<b>PER</b> person: <i>Axl_Rose, Sherlock_Holmes</i> <b>GPE</b> geo-political entity: <i>Paris, Japan</i> <b>GPO</b> *geo-political origin: <i>Parisian, French</i> <b>GEO</b> geographical location: <i>Alps, Nile</i> <b>ORG</b> organization: <i>IKEA, EU</i> <b>ART</b> artifact: <i>iOS_7</i> <b>HAP</b> happening: <i>Eurovision_2017</i> <b>UOM</b> unit of measurement: <i>meter, \$, %, degree_Celsius</i> <b>CTC</b> *contact information: <i>112, info@mail.com</i> <b>URL</b> URL: <i>http://pmb.let.rug.nl</i> <b>LIT</b> *literal use of names: <i>his name is John</i> <b>NTH</b> *other names: <i>table 1a, equation (1)</i>	<b>NAM</b> named entity
<b>COM</b> comparative	<b>EQU</b> equative: <i>as tall as John, whales are mammals</i> <b>MOR</b> comparative positive: <i>better, more</i> <b>LES</b> comparative negative: <i>less, worse</i>	<b>EXS</b> untensed simple: <i>to walk, is eaten, destruction</i> <b>ENS</b> present simple: <i>we walk, he walks</i> <b>EPS</b> past simple: <i>ate, went</i> <b>EXG</b> untensed progressive: <i>is running</i> <b>EXT</b> untensed perfect: <i>has eaten</i>	<b>EVE</b> events
<b>U</b> unna	<b>CON</b>	<b>NOW</b> present tense: <i>is skiing, do ski, has skied, now</i> <b>PST</b> past tense: <i>was baked, had gone, did go</i> <b>FUT</b> future tense: <i>will, shall</i> <b>PRG</b> *progressive: <i>has been being treated, aan.her<sup>NL</sup></i> <b>PFT</b> *perfect: <i>has been going/done</i>	<b>TNS</b> tense & aspect
<b>D</b> d	<b>ENS</b>	<b>DAT</b> *full date: <i>27.04.2017, 27/04/17</i> <b>DOM</b> day of month: <i>27th December</i> <b>YOC</b> year of century: <i>2017</i> <b>DOW</b> day of week: <i>Thursday</i> <b>MOY</b> month of year: <i>April</i> <b>DEC</b> decade: <i>80s, 1990s</i> <b>CLO</b> clocktime: <i>8:45_pm, 10_o'clock, noon</i>	<b>TIM</b> temporal entity
<b>DEF</b>	<b>GPO</b>		
<b>L</b> lo	<b>CON</b>		
	<b>NIL</b>		

DEF The  
CON yakuza  
ENS are  
DEF the  
GPO Japanese  
CON mafia  
NIL .



# Tagging

- The repository also has an example of training a semantic tagger
- Like POS tagging, but with a richer set of “semantic” tags
- Issue: the data comes with its own tokenization:
  - BERT: ['the', 'ya', '##zuka', 'are', 'the', 'japanese', 'mafia', '.']
- Need to get word-level representations out of BERT’s *subword* representations

<b>ANA</b> anaphoric	<b>PRO</b> anaphoric & deictic pronouns: <i>he, she, I, him</i> <b>DEF</b> definite: <i>the, lo<sup>IT</sup>, der<sup>DE</sup></i> <b>HAS</b> possessive pronoun: <i>my, her</i> <b>REF</b> reflexive & reciprocal pron.: <i>herself, each_other</i> <b>EMP</b> emphasizing pronouns: <i>himself</i>	<b>NOT</b> negation: <i>not, no, neither, without</i> <b>NEC</b> necessity: <i>must, should, have to</i> <b>POS</b> possibility: <i>might, could, perhaps, alleged, can</i>	<b>MOD</b> modality
<b>ACT</b> speech act	<b>GRE</b> greeting & parting: <i>hi, bye</i> <b>ITJ</b> interjections, exclamations: <i>alas, ah</i> <b>HES</b> hesitation: <i>err</i> <b>QUE</b> interrogative: <i>who, which, ?</i>	<b>SUB</b> subordinate relations: <i>that, while, because</i> <b>COO</b> coordinate relations: <i>so, {, }, {;}, and</i> <b>APP</b> appositional relations: <i>{, }, which, {(), {—}</i> <b>BUT</b> contrast: <i>but, yet</i>	<b>DSC</b> discourse
<b>ATT</b> attribute	<b>QUC</b> *concrete quantity: <i>two, six_million, twice</i> <b>QUV</b> *vague quantity: <i>millions, many, enough</i> <b>COL</b> *colour: <i>red, crimson, light_blue, chestnut_brown</i> <b>IST</b> intersective: <i>open, vegetarian, quickly</i> <b>SST</b> subsective: <i>skillful surgeon, tall kid</i> <b>PRI</b> privative: <i>former, fake</i> <b>DEG</b> *degree: <i>2 meters tall, 20 years old</i> <b>INT</b> intensifier: <i>very, much, too, rather</i> <b>REL</b> relation: <i>in, on, 's, of, after</i> <b>SCO</b> score: <i>3-0, grade A</i>	<b>PER</b> person: <i>Axl_Rose, Sherlock_Holmes</i> <b>GPE</b> geo-political entity: <i>Paris, Japan</i> <b>GPO</b> *geo-political origin: <i>Parisian, French</i> <b>GEO</b> geographical location: <i>Alps, Nile</i> <b>ORG</b> organization: <i>IKEA, EU</i> <b>ART</b> artifact: <i>iOS_7</i> <b>HAP</b> happening: <i>Eurovision_2017</i> <b>UOM</b> unit of measurement: <i>meter, \$, %, degree_Celsius</i> <b>CTC</b> *contact information: <i>112, info@mail.com</i> <b>URL</b> URL: <i>http://pmb.let.rug.nl</i> <b>LIT</b> *literal use of names: <i>his name is John</i> <b>NTH</b> *other names: <i>table 1a, equation (1)</i>	<b>NAM</b> named entity
<b>COM</b> comparative	<b>EQU</b> equative: <i>as tall as John, whales are mammals</i> <b>MOR</b> comparative positive: <i>better, more</i> <b>LES</b> comparative negative: <i>less, worse</i>	<b>EXS</b> untensed simple: <i>to walk, is eaten, destruction</i> <b>ENS</b> present simple: <i>we walk, he walks</i> <b>EPS</b> past simple: <i>ate, went</i> <b>EXG</b> untensed progressive: <i>is running</i> <b>EXT</b> untensed perfect: <i>has eaten</i>	<b>EVE</b> events
<b>U</b> unna	<b>CON</b>	<b>NOW</b> present tense: <i>is skiing, do ski, has skied, now</i> <b>PST</b> past tense: <i>was baked, had gone, did go</i> <b>FUT</b> future tense: <i>will, shall</i> <b>PRG</b> *progressive: <i>has been being treated, aan.her<sup>NL</sup></i> <b>PFT</b> *perfect: <i>has been going/done</i>	<b>TNS</b> tense & aspect
<b>D</b> d	<b>ENS</b>	<b>DAT</b> *full date: <i>27.04.2017, 27/04/17</i> <b>DOM</b> day of month: <i>27th December</i> <b>YOC</b> year of century: <i>2017</i> <b>DOW</b> day of week: <i>Thursday</i> <b>MOY</b> month of year: <i>April</i> <b>DEC</b> decade: <i>80s, 1990s</i> <b>CLO</b> clocktime: <i>8:45_pm, 10_o'clock, noon</i>	<b>TIM</b> temporal entity
<b>DEF</b>	<b>GPO</b>		
<b>CON</b>	<b>NIL</b>		

# Tagging: Modeling

- Used to be complicated, BUT:
- They've added a [PretrainedMismatchedTransformerEmbedder](#) (and a corresponding [PretrainedMismatchedTransformerIndexer](#) for tokens)
- Handles all of the mis-alignment between dataset tokens and model tokens for you!
- How to pool subwords—>words:
  - ``sub_token_mode`` kwarg: default = avg, but can do first/last, etc

# Tagging: Modeling

```
@Model.register("subword_word_tagger")
class SubwordWordTagger(Model):

    # TODO: document!

    def __init__(
        self,
        embedder: TextFieldEmbedder,
        vocab: Vocabulary = None,
        freeze_encoder: bool = True,
    ):
        super().__init__(vocab)

        self._embedder = embedder
        self._freeze_encoder = freeze_encoder
        # turn off gradients if don't want to fine tune encoder
        for parameter in self._embedder.parameters():
            parameter.requires_grad = not self._freeze_encoder

        self.classifier = TimeDistributed(
            torch.nn.Linear(
                in_features=embedder.get_output_dim(),
                out_features=vocab.get_vocab_size("labels"),
            )
        )

        self.accuracy = CategoricalAccuracy()
```

```
def forward(
    self, sentence: Dict[str, torch.Tensor], labels: torch.Tensor = None
) -> Dict[str, torch.Tensor]:

    # (batch_size, max_seq_len, embedding_dim)
    embeddings = self._embedder(sentence)

    # get mask to keep track of which tokens are padding
    # prevent them from contributing to loss
    # (batch_size, max_word_seq_len)
    word_mask = get_text_field_mask(sentence)

    # (batch_size, max_word_seq_len, num_labels)
    logits = self.classifier(embeddings)

    outputs = {"logits": logits}

    if labels is not None:
        self.accuracy(logits, labels, word_mask)
        outputs["loss"] = sequence_cross_entropy_with_logits(
            logits, labels, word_mask
        )

    return outputs
```



# On These Libraries

- If you're using transformer-based LMs, I strongly recommend HuggingFace
- On the other hand, it's possible that learning AllenNLP's abstractions may cost you more time than it saves in the short term
- As always, try and use the best tool for the job at hand
- One more that makes fine-tuning and/or diagnostic classification easy:
  - [jiant](#)

# Other tools for experiment management

- Disclaimer: I've never used them!
  - Might be over-kill in the short term
- Guild (entirely local): <https://guild.ai/>
- CodaLab: <https://codalab.org/>
- Weights and Biases: <https://www.wandb.com/>
- Neptune: <https://neptune.ai/>

# Using GPUs on Patas

# Setting up local environment

- Three GPU nodes:
  - 2xTesla P40
  - 8xTesla M10
  - 2xQuadro 8000
- For info on setting up your local environment to use these nodes in a fairly painless way:
  - <https://www.shane.st/teaching/575/spr22/patas-gpu.pdf>



# Condor job file for patas

```
executable = run_exp_gpu.sh
getenv = True
error = exp.error
log = exp.log
notification = always
transfer_executable = false
request_memory = 8*1024
request_GPUs = 1
+Research = True
Queue
```

# Example executable

```
#!/bin/sh
conda activate my-project

allennlp train tagging_experiment.jsonnet --serialization-dir test \
  --include-package tagging \
  --overrides '{"trainer": {"cuda_device": 1}}'
```