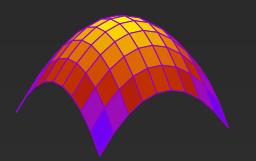


State-of-the-Art Transformer Pipelines with spaCy

Daniël de Kok



Biaffine

Madeesh Kannan

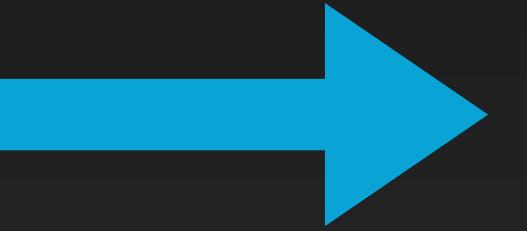


Explosion

10th November 2023, aiGrunn

What is spaCy?

Text documents



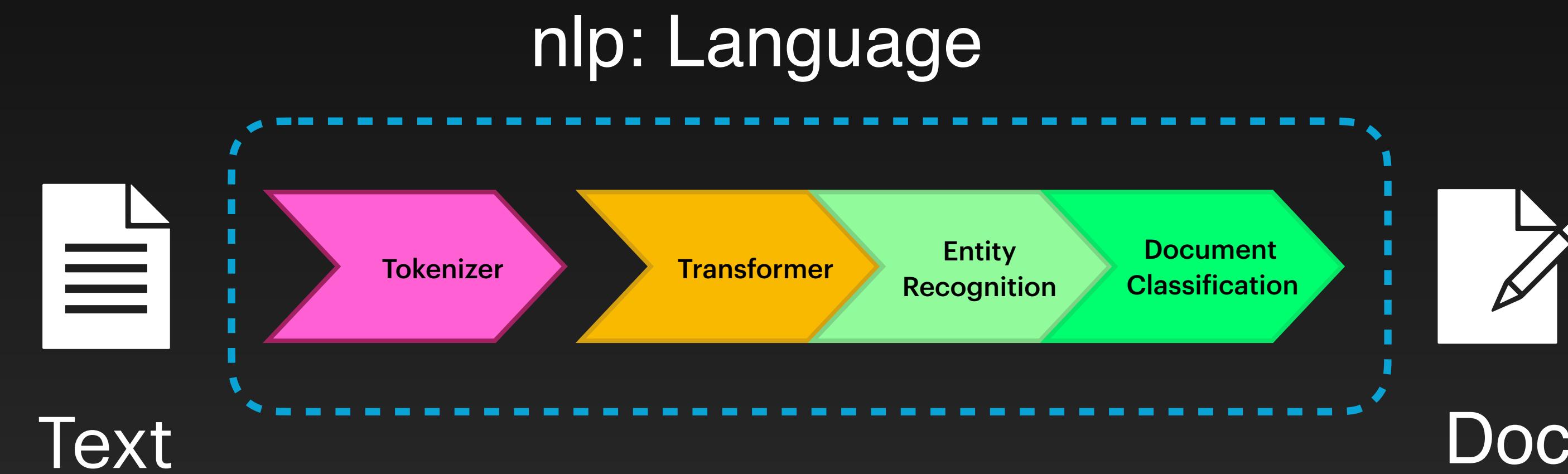
Annotated documents



spaCy

```
$ pip install spacy
```

spaCy pipeline



Document container



On April 25, Twitter's board of directors

unanimously accepted Musk's buyout offer.

Tokenization



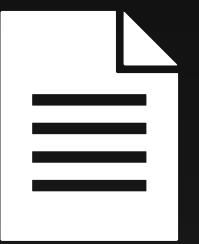
doc: Doc

On April 25, Twitter's board of directors
unanimously accepted Musk's buyout offer.

`class Token:`
`text: str`
`idx: int`

doc[9]: Token

Lemmatization



doc: Doc

on

On April 25, Twitter's board of directors

accept

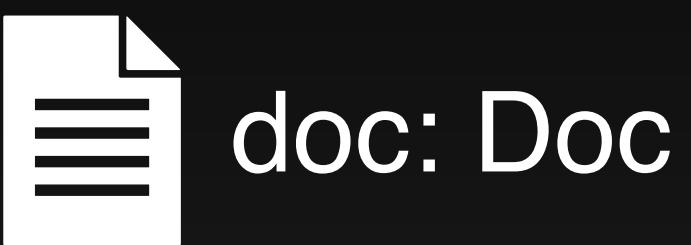
unanimously accepted Musk's buyout offer.

director

doc[10].lemma_: str

class Token:
text: str
lemma_: str
idx: int

Span classification: entity recognition



doc: Doc

DATE

ORG

On April 25, Twitter's board of directors

PERSON

unanimously accepted Musk's buyout offer.

doc.ents[2]: Span

`class Span:`

`text: str`

`label_: str`

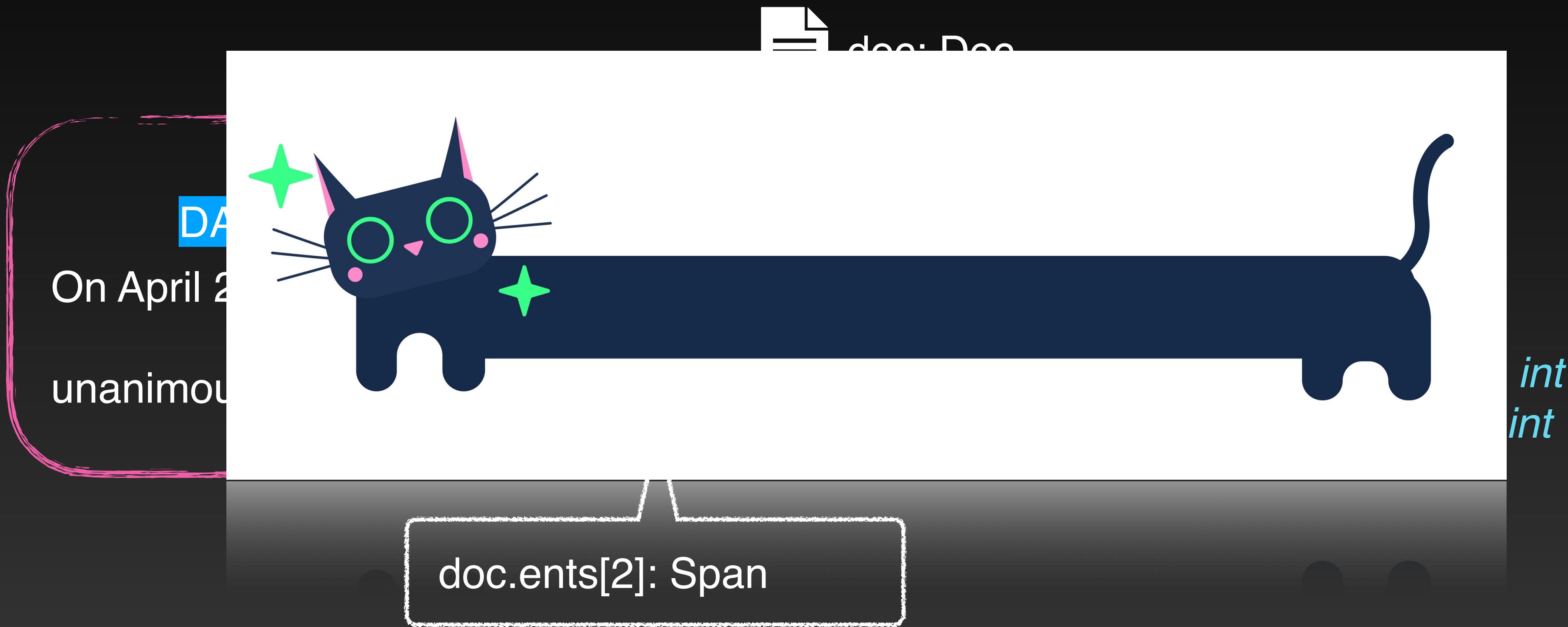
`start: int`

`end: int`

`start_char: int`

`end_char: int`

Span classification: entity recognition



Document classification



doc: Doc

On April 25, Twitter's board of directors
unanimously accepted Musk's buyout offer.

Newswire

doc.cats: Dict[str, float]

How transformers are used

Hidden Representations



Station to Station is the tenth studio album by David Bowie .

SpanCat

Lemmatizer

How transformers are used



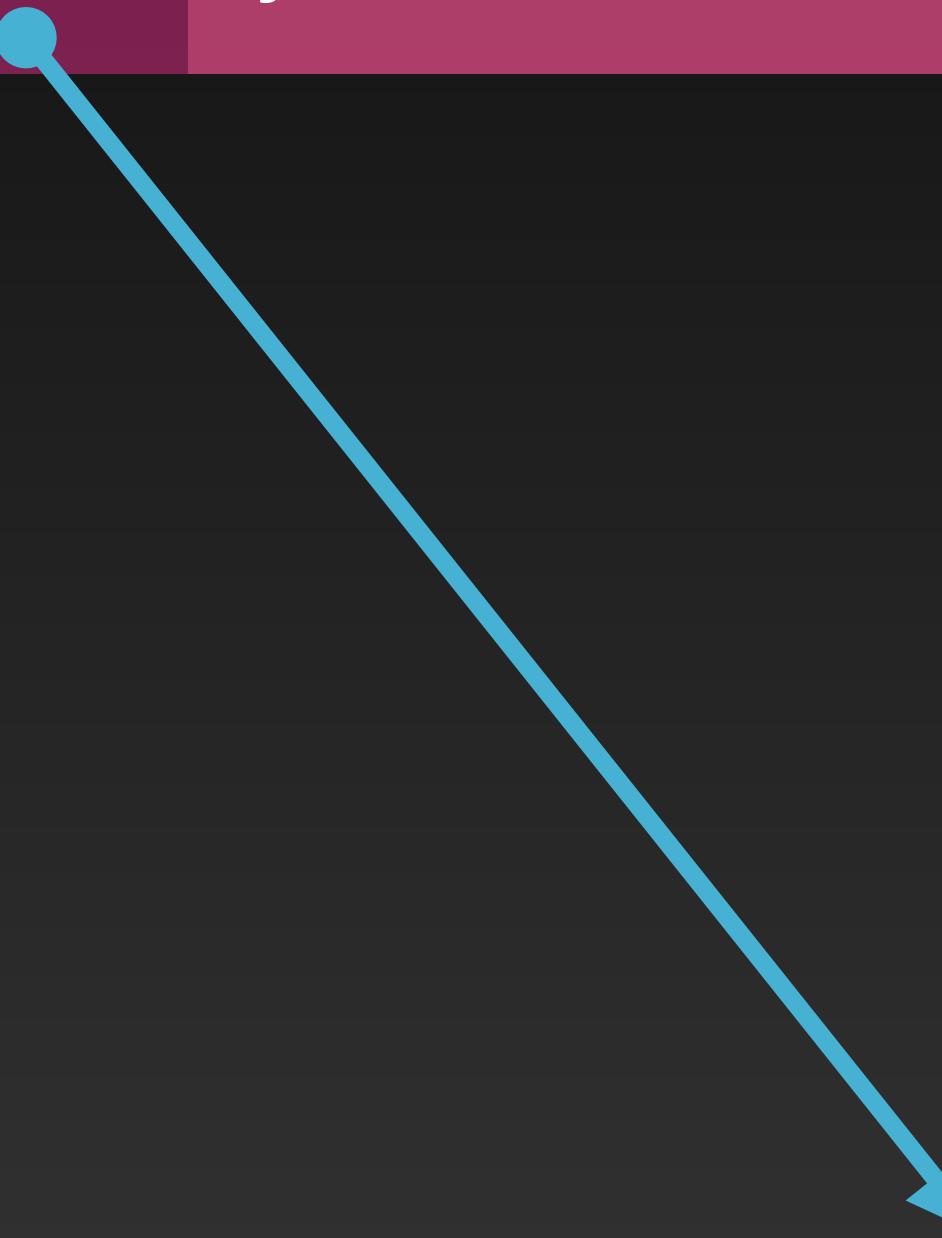
Hidden Representations

Station to Station is the tenth studio album by David Bowie .

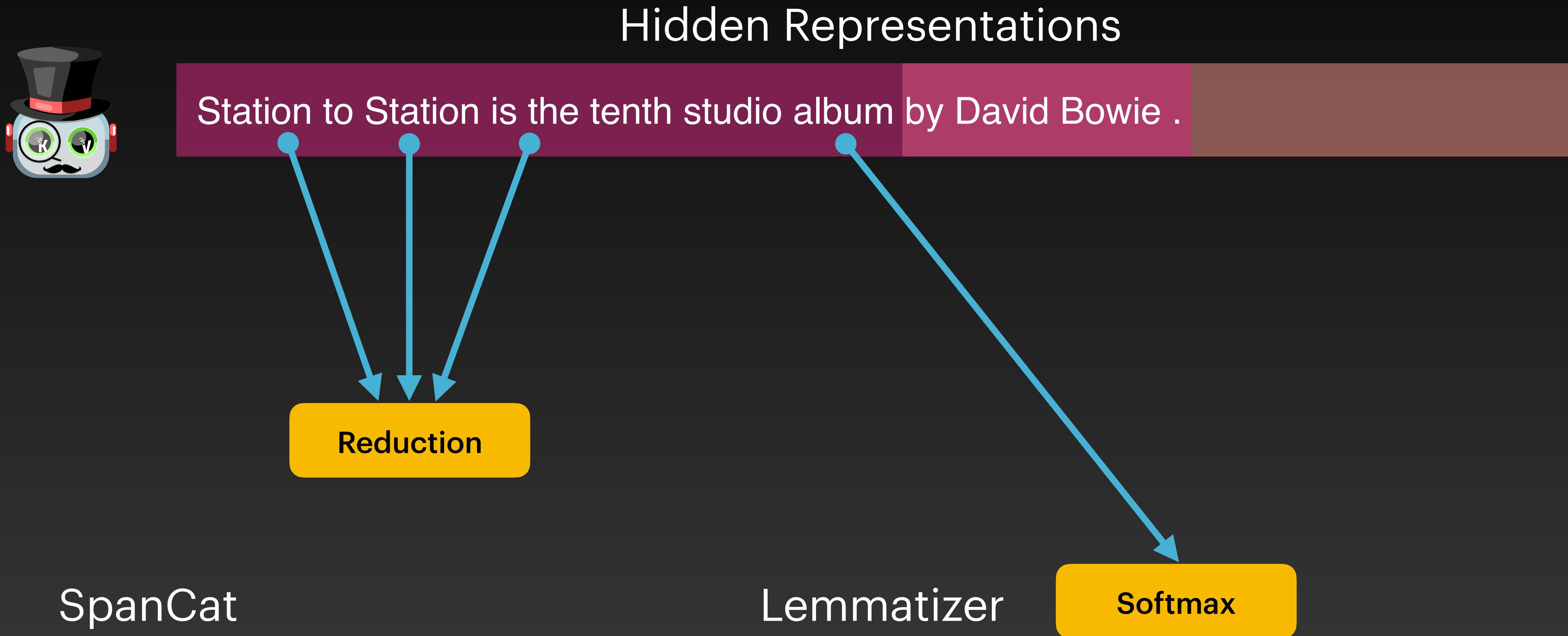
SpanCat

Lemmatizer

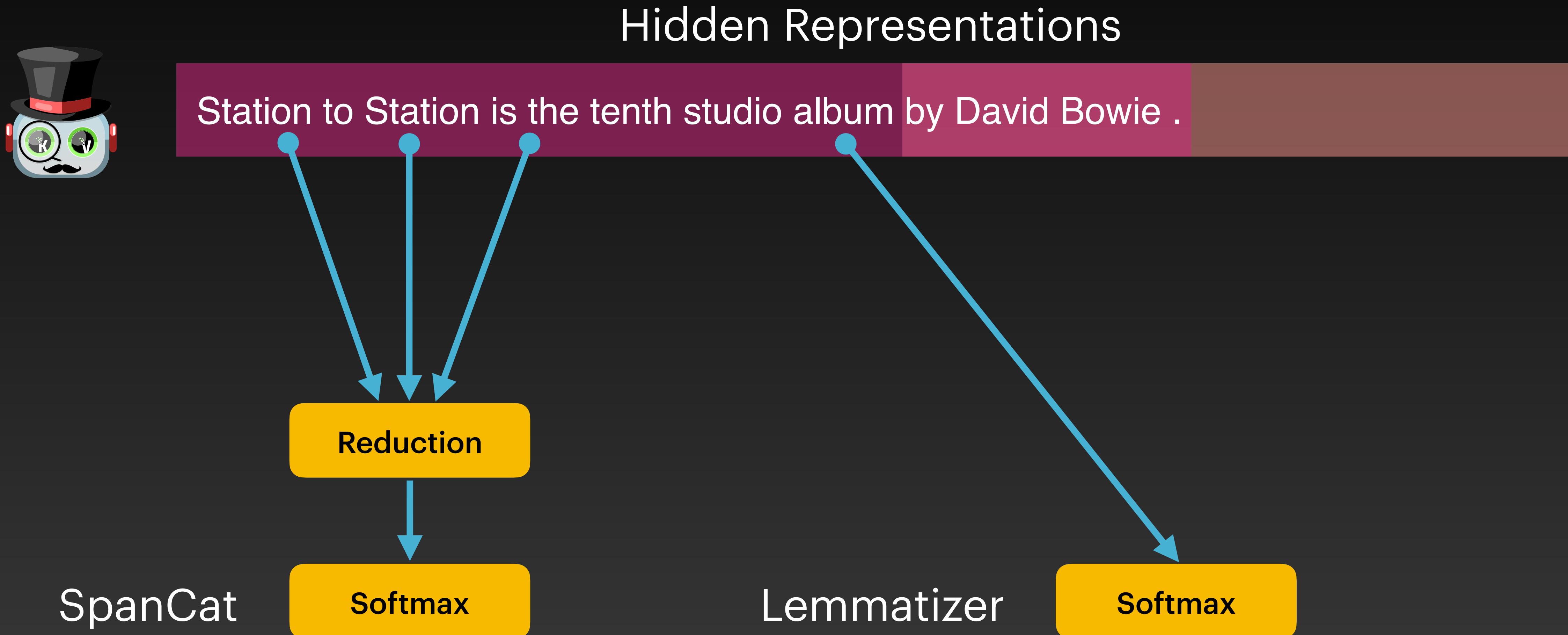
Softmax



How transformers are used



How transformers are used



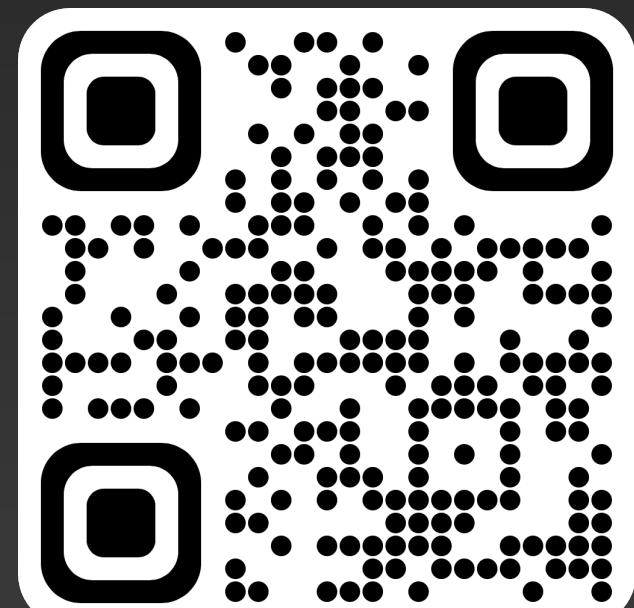
Inference with Transformer Pipelines

Pretrained transformer models

Training data
en_core_web_trf
Language

```
$ spacy download en_core_web_trf
```

```
nlp = spacy.load("en_core_web_trf")
```



Small webservice

- Small FastAPI REST web service with /annotate endpoint.
- The user provides texts in JSON format.
- Process the texts with the en_core_web_trf pipeline.
- Return:
 - Tokens
 - Lemmas
 - Entities

Endpoint

```
en_core_web_trf = spacy.load("en_core_web_trf")

@app.post("/annotate")
def annotate(docs: List[str]):
    annotated_docs = en_core_web_trf.pipe(docs)
    return [convert_spacy_doc(doc) for doc in annotated_docs]
```

Data model

```
@dataclass  
class Token:  
    text: str  
    lemma: str
```

```
@dataclass  
class Span:  
    text: str  
    label: str  
    span: Tuple[int, int]
```

```
@dataclass  
class Annotations:  
    tokens: List[Token]  
    entities: List[Span]
```

Data conversion

```
def convert_spacy_doc(doc: Doc):
    tokens = [Token(text=token.text, lemma=token.lemma_)
              for token in doc]
    entities = [
        Span(text=ent.text,
              label=ent.label_,
              span=(ent.start, ent.end))
    ]
    for ent in doc.ents
    ]
    return Annotations(tokens=tokens, entities=entities)
```

Example run

```
% curl \  
-X POST \  
-H "Content-Type: application/json" \  
-d '["The Glove80 was released by MoErgo in 2023, bringing joy to  
many typists."'] \  
http://localhost:8000/annotate | jq '.[].entities'
```

```
[  
 {  
   "text": "Glove80",  
   "label": "PRODUCT",  
   "span": [1, 2]  
 },  
 {  
   "text": "MoErgo",  
   "label": "ORG",  
   "span": [5, 6]  
 },  
 {  
   "text": "2023",  
   "label": "DATE",  
   "span": [7, 8]  
 }]
```

Training Transformer Pipelines

Training data

- Groningen Meaning Bank (GMB):
 - Developed by the University of Groningen.
 - >10K annotated English texts.
 - Genres: Predominantly newspaper text.
 - Public domain.
 - Annotations: Named Entities, Lemmas, Part-of-Speech



Pipeline configuration

- Full description of the pipeline.
- Structured, pluggable sections.
- Reproducibility.
- Validation.



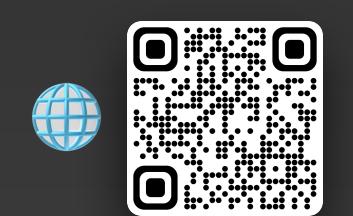
```
1 [paths]
2 train = null
3 dev = null
4
5 [system]
6 gpu_allocator = "pytorch"
7 seed = 1
8
9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```



Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

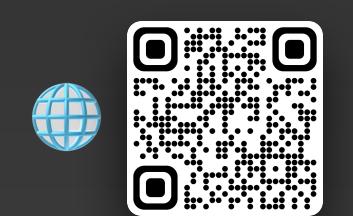


Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

[**paths**] - Paths to data and other assets.

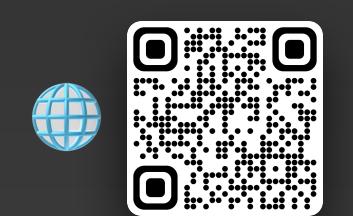


Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

[system] - Settings related to system and hardware.

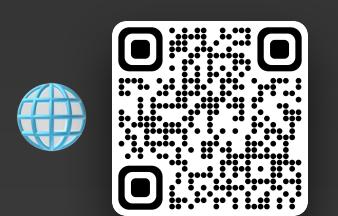


Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

[nlp] - Pipeline components, language and tokenizer.

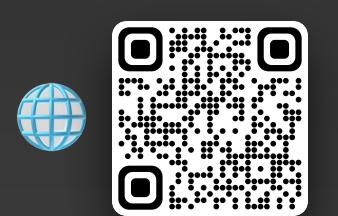


Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

[components] - Definitions of pipeline components.

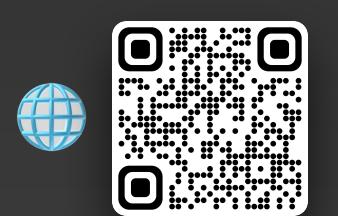


Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

[corpora] - Readers for the training data.

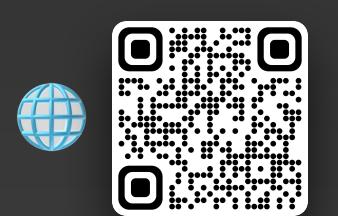


Pipeline configuration

Sections

```
 1 [paths]
 2 train = null
 3 dev = null
 4
 5 [system]
 6 gpu_allocator = "pytorch"
 7 seed = 1
 8
 9 [nlp]
10 lang = "en"
11 pipeline = ["transformer", "tagger", "lemmatizer", "ner"]
12 disabled = []
13 batch_size = 64
14 tokenizer = {"@tokenizers": "spacy.Tokenizer.v1"}
15
16 [components]
17
18 [components.lemmatizer]
19 factory = "trainable_lemmatizer"
20 min_tree_freq = 3
21 overwrite = false
22 scorer = {"@scorers": "spacy.lemmatizer_scorer.v1"}
23 top_k = 1
24
25 [corpora]
26
27 [corpora.train]
28 @readers = "spacy.Corporus.v1"
29 path = ${paths.train}
30 gold_preproc = false
31 max_length = 0
32 limit = 0
33
34 [training]
35 train_corpus = "corpora.train"
36 dev_corpus = "corpora.dev"
37 seed = ${system:seed}
38 gpu_allocator = ${system:gpu_allocator}
39 dropout = 0.1
40 accumulate_gradient = 3
41 patience = 5000
42 max_epochs = 5
43 max_steps = 10000
44 eval_frequency = 250
```

[**training**] - Training parameters.



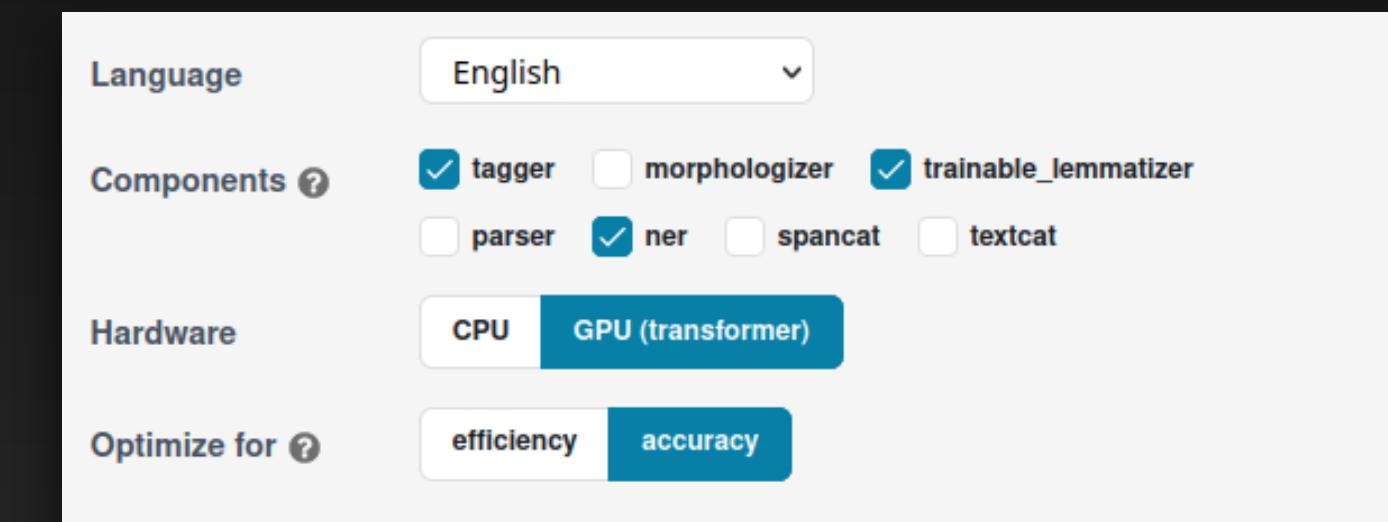
Pipeline configuration

Command-line:

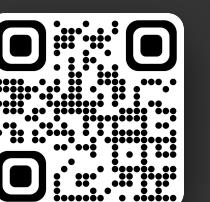
```
$ spacy init config base.cfg
```

```
$ spacy init fill-config base.cfg  
filled.cfg
```

Quickstart widget on spacy.io

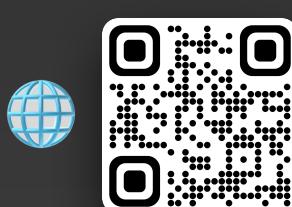
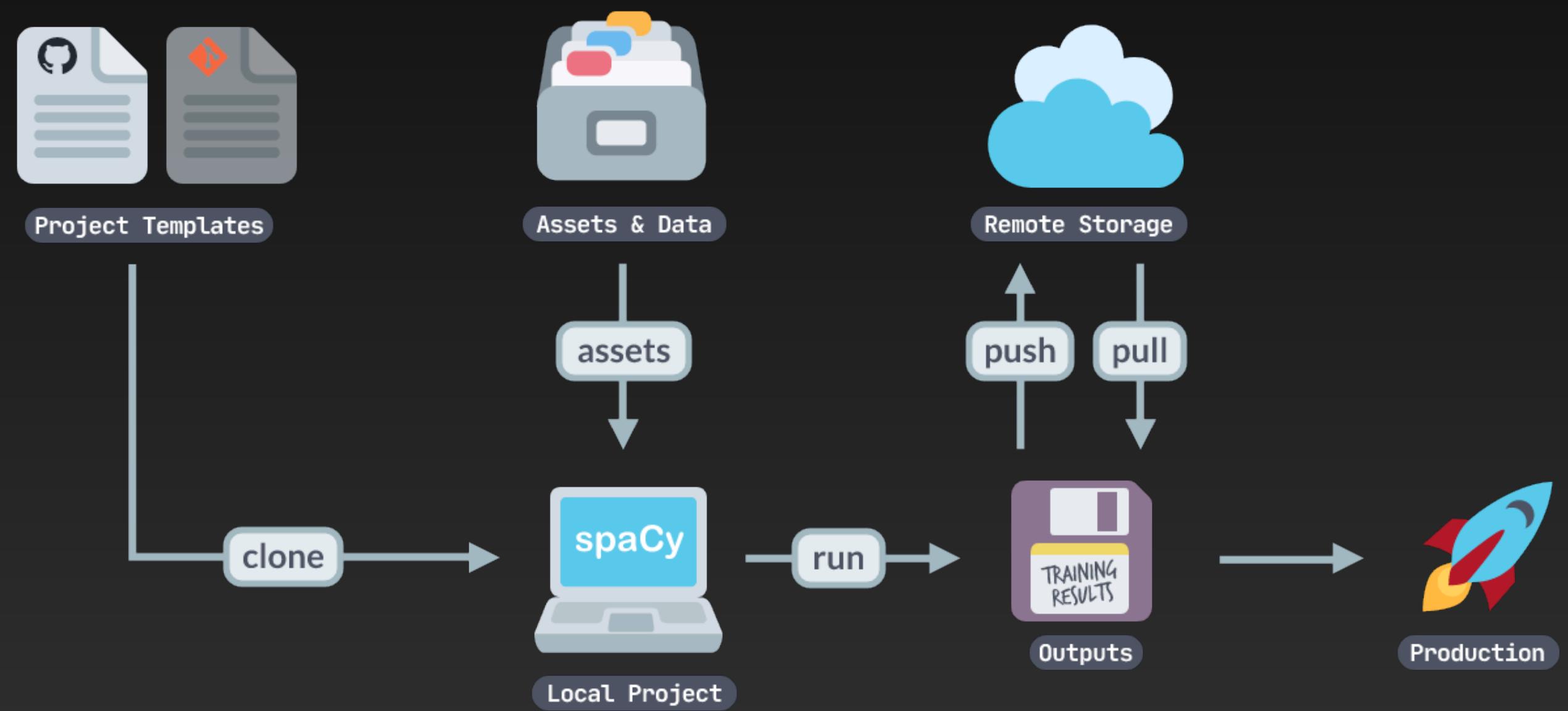


```
# This is an auto-generated partial config. To use it with 'spacy  
train'  
# you can run spacy init fill-config to auto-fill all default  
settings:  
# python -m spacy init fill-config ./base_config.cfg ./config.cfg  
[paths]  
train = null  
dev = null  
vectors = null  
[system]  
gpu_allocator = "pytorch"  
  
[nlp]  
lang = "en"  
pipeline = ["transformer", "tagger", "ner", "trainable_lemmatizer"]  
batch_size = 128  
  
[components]
```



Projects

- End-to-end spaCy workflows.
 - Assets.
 - Training.
 - Evaluation.
 - Deployment.

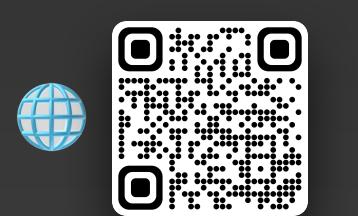


Projects

Definition



```
1 title: "spaCy Transformer POS/Lemmatizer/NER pipeline"
2
3 # Variables can be referenced across the project.yml using ${vars.var_name}
4 vars:
5   seed: 1
6   splits: "80|10|10"
7
8
9 directories: ["assets", "training", "configs", "metrics", "corpus"]
10
11
12 assets:
13   - dest: "assets/gmb.zip"
14     url: "https://gmb.let.rug.nl/releases/gmb-2.2.0.zip"
15     checksum: "7b12fd710826d5a4963e57a43100f6c6"
16
17 commands:
18   - name: corpus
19     help: "Convert the data to spaCy's format"
20     script:
21       - 'unzip assets/gmb.zip "*.tags" -d assets/'
22       - >-
23         python scripts/conv.py assets/gmb-2.2.0 corpus
24         --seed ${vars.seed}
25         --splits '${vars.splits}'
26     deps:
27       - "assets/gmb.zip"
28     outputs:
29       - "corpus/train.spacy"
30       - "corpus/dev.spacy"
31       - "corpus/test.spacy"
32
33 workflows:
34   all:
35     - corpus
36     - train
37     - evaluate
```



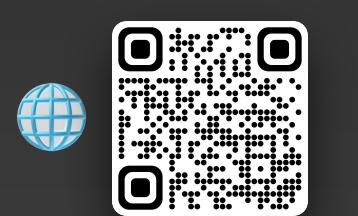
Projects

Definition



```
1 title: "spaCy Transformer POS/Lemmatizer/NER pipeline"
2
3 # Variables can be referenced across the project.yml using ${vars.var_name}
4 vars:
5   seed: 1
6   splits: "80|10|10"
7
8
9 directories: ["assets", "training", "configs", "metrics", "corpus"]
10
11
12 assets:
13   - dest: "assets/gmb.zip"
14     url: "https://gmb.let.rug.nl/releases/gmb-2.2.0.zip"
15     checksum: "7b12fd710826d5a4963e57a43100f6c6"
16
17 commands:
18   - name: corpus
19     help: "Convert the data to spaCy's format"
20     script:
21       - 'unzip assets/gmb.zip "*.tags" -d assets/'
22       - >-
23         python scripts/conv.py assets/gmb-2.2.0 corpus
24         --seed ${vars.seed}
25         --splits '${vars.splits}'
26     deps:
27       - "assets/gmb.zip"
28     outputs:
29       - "corpus/train.spacy"
30       - "corpus/dev.spacy"
31       - "corpus/test.spacy"
32
33 workflows:
34   all:
35     - corpus
36     - train
37     - evaluate
```

vars - Configurable settings.



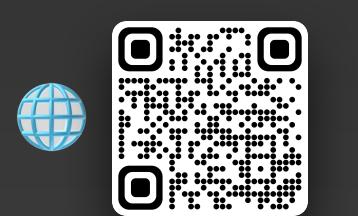
Projects

Definition



```
1 title: "spaCy Transformer POS/Lemmatizer/NER pipeline"
2
3 # Variables can be referenced across the project.yml using ${vars.var_name}
4 vars:
5   seed: 1
6   splits: "80|10|10"
7
8
9 directories: ["assets", "training", "configs", "metrics", "corpus"]
10
11
12 assets:
13   - dest: "assets/gmb.zip"
14     url: "https://gmb.let.rug.nl/releases/gmb-2.2.0.zip"
15     checksum: "7b12fd710826d5a4963e57a43100f6c6"
16
17 commands:
18   - name: corpus
19     help: "Convert the data to spaCy's format"
20     script:
21       - 'unzip assets/gmb.zip "*.tags" -d assets/'
22       - >-
23         python scripts/conv.py assets/gmb-2.2.0 corpus
24         --seed ${vars.seed}
25         --splits '${vars.splits}'
26     deps:
27       - "assets/gmb.zip"
28     outputs:
29       - "corpus/train.spacy"
30       - "corpus/dev.spacy"
31       - "corpus/test.spacy"
32
33 workflows:
34   all:
35     - corpus
36     - train
37     - evaluate
```

directories - Directories for outputs, assets, etc.



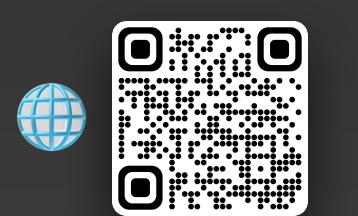
Projects

Definition



```
1 title: "spaCy Transformer POS/Lemmatizer/NER pipeline"
2
3 # Variables can be referenced across the project.yml using ${vars.var_name}
4 vars:
5   seed: 1
6   splits: "80|10|10"
7
8
9 directories: ["assets", "training", "configs", "metrics", "corpus"]
10
11
12 assets:
13   - dest: "assets/gmb.zip"
14     url: "https://gmb.let.rug.nl/releases/gmb-2.2.0.zip"
15     checksum: "7b12fd710826d5a4963e57a43100f6c6"
16
17 commands:
18   - name: corpus
19     help: "Convert the data to spaCy's format"
20     script:
21       - 'unzip assets/gmb.zip "*.tags" -d assets/'
22       - >-
23         python scripts/conv.py assets/gmb-2.2.0 corpus
24         --seed ${vars.seed}
25         --splits '${vars.splits}'
26     deps:
27       - "assets/gmb.zip"
28     outputs:
29       - "corpus/train.spacy"
30       - "corpus/dev.spacy"
31       - "corpus/test.spacy"
32
33 workflows:
34   all:
35     - corpus
36     - train
37     - evaluate
```

assets - List of asset definitions.



Projects

Definition



```
1 title: "spaCy Transformer POS/Lemmatizer/NER pipeline"
2
3 # Variables can be referenced across the project.yml using ${vars.var_name}
4 vars:
5   seed: 1
6   splits: "80|10|10"
7
8
9 directories: ["assets", "training", "configs", "metrics", "corpus"]
10
11
12 assets:
13   - dest: "assets/gmb.zip"
14     url: "https://gmb.let.rug.nl/releases/gmb-2.2.0.zip"
15     checksum: "7b12fd710826d5a4963e57a43100f6c6"
16
17 commands:
18   - name: corpus
19     help: "Convert the data to spaCy's format"
20     script:
21       - 'unzip assets/gmb.zip "*.tags" -d assets/'
22       - >-
23         python scripts/conv.py assets/gmb-2.2.0 corpus
24         --seed ${vars.seed}
25         --splits '${vars.splits}'
26     deps:
27       - "assets/gmb.zip"
28     outputs:
29       - "corpus/train.spacy"
30       - "corpus/dev.spacy"
31       - "corpus/test.spacy"
32
33 workflows:
34   all:
35     - corpus
36     - train
37     - evaluate
```

commands - List of named commands, their dependencies and outputs.



Projects

Definition



```
1 title: "spaCy Transformer POS/Lemmatizer/NER pipeline"
2
3 # Variables can be referenced across the project.yml using ${vars.var_name}
4 vars:
5   seed: 1
6   splits: "80|10|10"
7
8
9 directories: ["assets", "training", "configs", "metrics", "corpus"]
10
11
12 assets:
13   - dest: "assets/gmb.zip"
14     url: "https://gmb.let.rug.nl/releases/gmb-2.2.0.zip"
15     checksum: "7b12fd710826d5a4963e57a43100f6c6"
16
17 commands:
18   - name: corpus
19     help: "Convert the data to spaCy's format"
20     script:
21       - 'unzip assets/gmb.zip "*.tags" -d assets/'
22       - >-
23         python scripts/conv.py assets/gmb-2.2.0 corpus
24         --seed ${vars.seed}
25         --splits '${vars.splits}'
26     deps:
27       - "assets/gmb.zip"
28     outputs:
29       - "corpus/train.spacy"
30       - "corpus/dev.spacy"
31       - "corpus/test.spacy"
32
33 workflows:
34   all:
35     - corpus
36     - train
37     - evaluate
```

workflows - Lists of commands to execute in-order.



Let's train a **spaCy** pipeline with:
Lemmatiser, **Named Entity Recogniser** and **Part-of-Speech Tagger**

Training Project

Structure:



Commands:

- 📚 corpus
- 💪 train
- ✅ evaluate

Training

Step 1: Download assets

Training

Step 1: Download assets



```
$ cd project  
$ spacy project assets  
  
i Fetching 1 asset(s)  
✓ Downloaded asset  
/Users/mkannan/dev/aiGrunn-2023/project/assets/gmb.zip
```

```
$ cd project  
$ spacy project assets
```

Training

Step 2: Convert corpus into training data

Training

Step 2: Convert corpus into training data



```
$ spacy project run corpus

=====
corpus =====
Running command: unzip assets/gmb.zip '*.tags' -d assets/
Archive: assets/gmb.zip
  inflating: assets/gmb-2.2.0/data/p76/d0310/en.tags
  inflating: assets/gmb-2.2.0/data/p76/d0689/en.tags
  inflating: assets/gmb-2.2.0/data/p76/d0687/en.tags
...
Running command: /Users/mkannan/opt/miniconda3/bin/python scripts/gmb_converter.py assets/gmb-2.2.0 corpus
--seed 1 --splits '80|10|10' --n-docs 5000

Input directory: assets/gmb-2.2.0
Output directory: corpus
Splits: (80, 10, 10)
Number of documents to convert: 5000
Coalesce consecutive named entities: True
Processing documents...
100% |██████████| 10000/10000 [00:08<00:00, 1186.21it/s]
Converting 5000 documents...
Wrote 4000 documents to the training set
Wrote 500 documents to the development set
Wrote 500 documents to the test set
```

\$ spacy project run corpus

Training

Step 3: Train the pipeline

Training

Step 3: Train the pipeline

```
● ● ●

$ spacy project run train

=====
train
=====
Running command: /Users/mkannan/opt/miniconda3/bin/python -m spacy train configs/default.cfg -o training/ --
gpu-id 0 --paths.train corpus/train.spacy --paths.dev corpus/dev.spacy
i Saving to output directory: training
i Using GPU: 0

=====
Initializing pipeline
=====
✓ Initialized pipeline

=====
Training pipeline
=====
i Pipeline: ['transformer', 'tagger', 'lemmatizer', 'ner']
i Initial learn rate: 0.0
E   #      LOSS TRANS...  LOSS TAGGER  LOSS LEMMA  LOSS NER  TAG_ACC  LEMMA_ACC  ENTS_F  SPEED  SCORE
---  ---  -----  -----  -----  -----  -----  -----  -----  -----  -----
0    0     11019.05    700.33    715.70    516.79    0.10    64.94    0.63    35326.51    0.30
0    250   541568.05   409296.48   403073.66   101223.78   54.00    66.10    76.57    38424.30    0.70
1    500   35135.77    99152.09   115594.67   22673.97   97.38    89.71    80.39    36970.75    0.86
2    750   23021.88    14806.34   37780.36   19576.18   98.38    96.88    83.10    35357.75    0.91
3   1000   39303.60    10942.97   15616.03   17299.31   98.42    98.20    83.85    41808.22    0.92
4   1250   31570.73    9091.26   9232.81   16120.65   98.51    98.64    84.20    41526.73    0.92
Epoch 5: 96%|██████████| 241/250 [01:45<00:04, 2.21it/s]
✓ Saved pipeline to output directory
training/model-last
```

\$ spacy project run train

Training

Step 4: Evaluate on held-out data

Training

Step 4: Evaluate on held-out data

```
● ● ●  
$ spacy project run evaluate  
  
===== evaluate =====  
Running command: /Users/mkannan/opt/miniconda3/bin/python -m spacy evaluate  
./training/model-best ./corpus/test.spacy --output ./metrics/default.json --gpu-id 0  
i Using GPU: 0  
  
===== Results =====  
TOK      -  
TAG     98.39  
LEMMA   98.66  
NER P    84.69  
NER R    84.32  
NER F    84.50  
SPEED   21740  
  
===== NER (per type) =====  
          P      R      F  
GEO-NAM  87.16  87.51  87.34  
TIM-DOW  98.78  100.00 99.39  
PER-NAM  71.52  78.80  74.98  
TIM-DAT  75.17  65.22  69.84  
TIM-MOY  99.01  100.00 99.50  
GPE-NAM  96.38  92.65  94.48  
TIM-CLO  90.00  45.00  60.00  
ART-NAM  0.00   0.00   0.00  
ORG-NAM  67.94  70.99  69.43  
PER-FAM  87.16  86.96  87.06  
PER-GIV  87.07  73.72  79.84  
  
✓ Saved results to metrics/default.json
```

\$ spacy project run evaluate

Using the trained model



```
import spacy

# Load the pipeline we just trained.
nlp = spacy.load("training/model-last")

# Process whole documents
text = ("When Sebastian Thrun started working on self-driving cars at "
        "Google in 2007, few people outside of the company took him "
        "seriously. "I can tell you very senior CEOs of major American "
        "car companies would shake my hand and turn away because I wasn't "
        "worth talking to," said Thrun, in an interview with Recode earlier "
        "this week.")
doc = nlp.pipe([text])

# Analyze syntax
print("Noun phrases:", [chunk.text for chunk in doc.noun_chunks])
print("Verbs:", [token.lemma_ for token in doc if token.pos_ == "VERB"])

# Find named entities, phrases and concepts
for entity in doc.ents:
    print(entity.text, entity.label_)
```

★ spaCy - spacy.io

★ Explosion - explosion.ai

▲ Biaffine - biaffine.ai

✉ Daniël - danieldk.eu

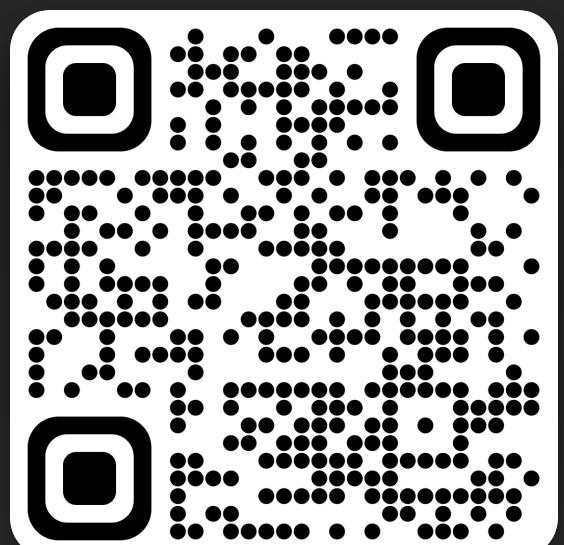
in Madeesh - [@m-kannan](https://www.linkedin.com/in/m-kannan)

Play around with the project!



github.com/explosion/aiGrunn-2023

Curated Transformers



★ spaCy - spacy.io

★ Explosion - explosion.ai

▲ Biaffine - biaffine.ai

✉ Daniël - danieldk.eu

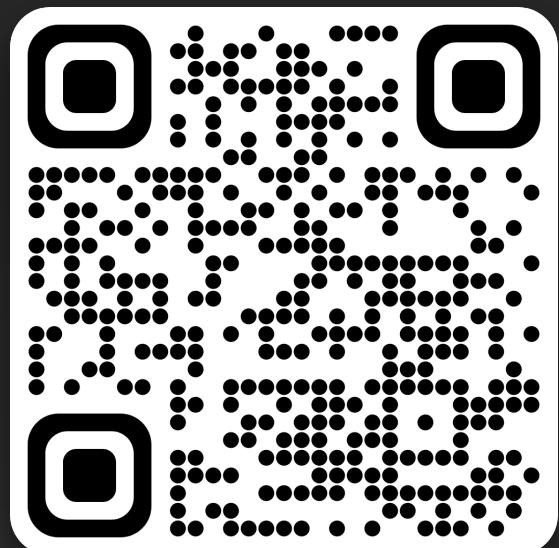
in Madeesh - [@m-kannan](https://www.linkedin.com/in/m-kannan)

Play around with the project!



github.com/explosion/aiGru...-2023

Curated Transformers A small, stylized cartoon character of a robot head with large green eyes, a mustache, and a top hat, positioned next to the text.

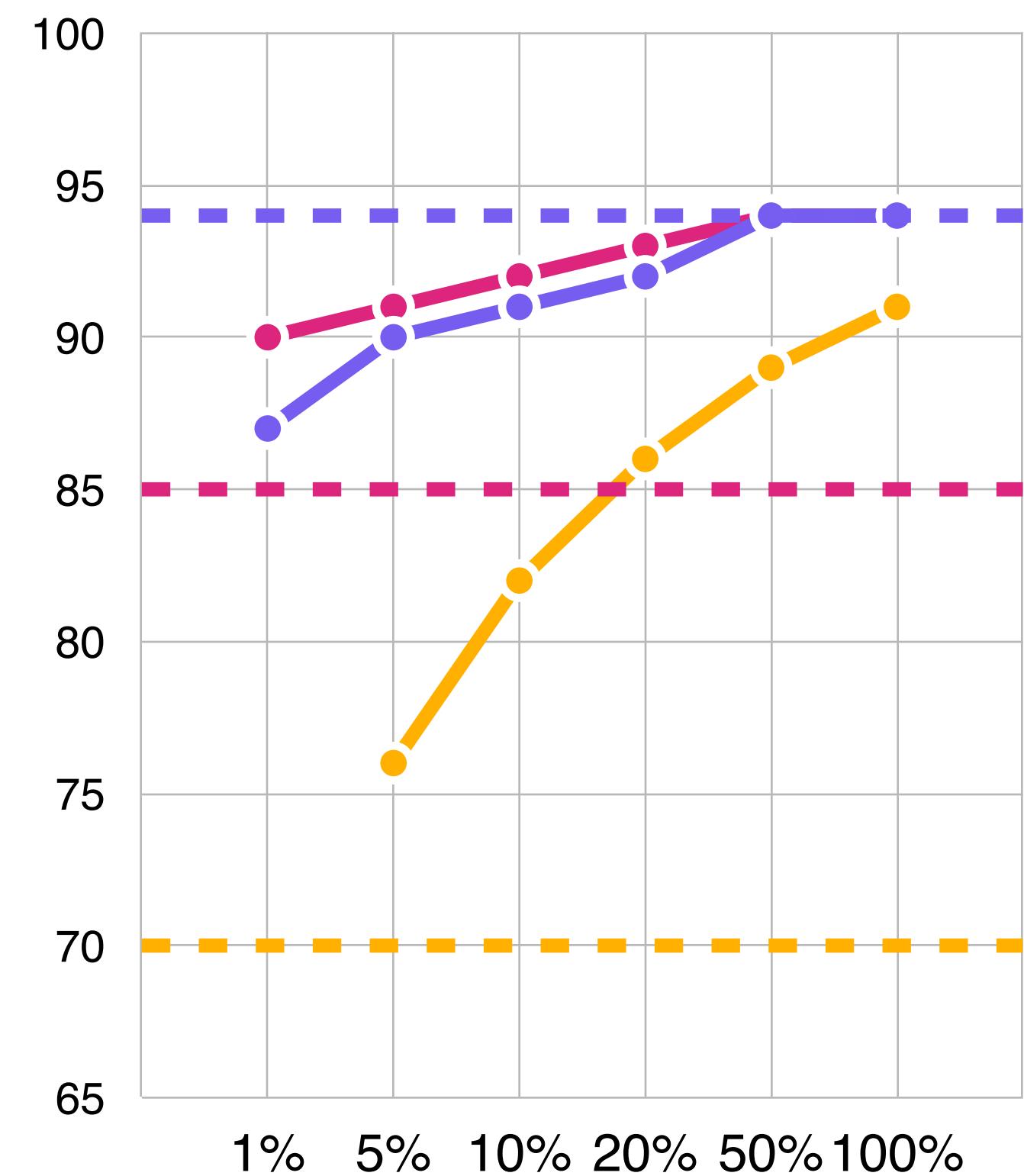


\$ pip install curated-transformers

EXPLO
NOISE

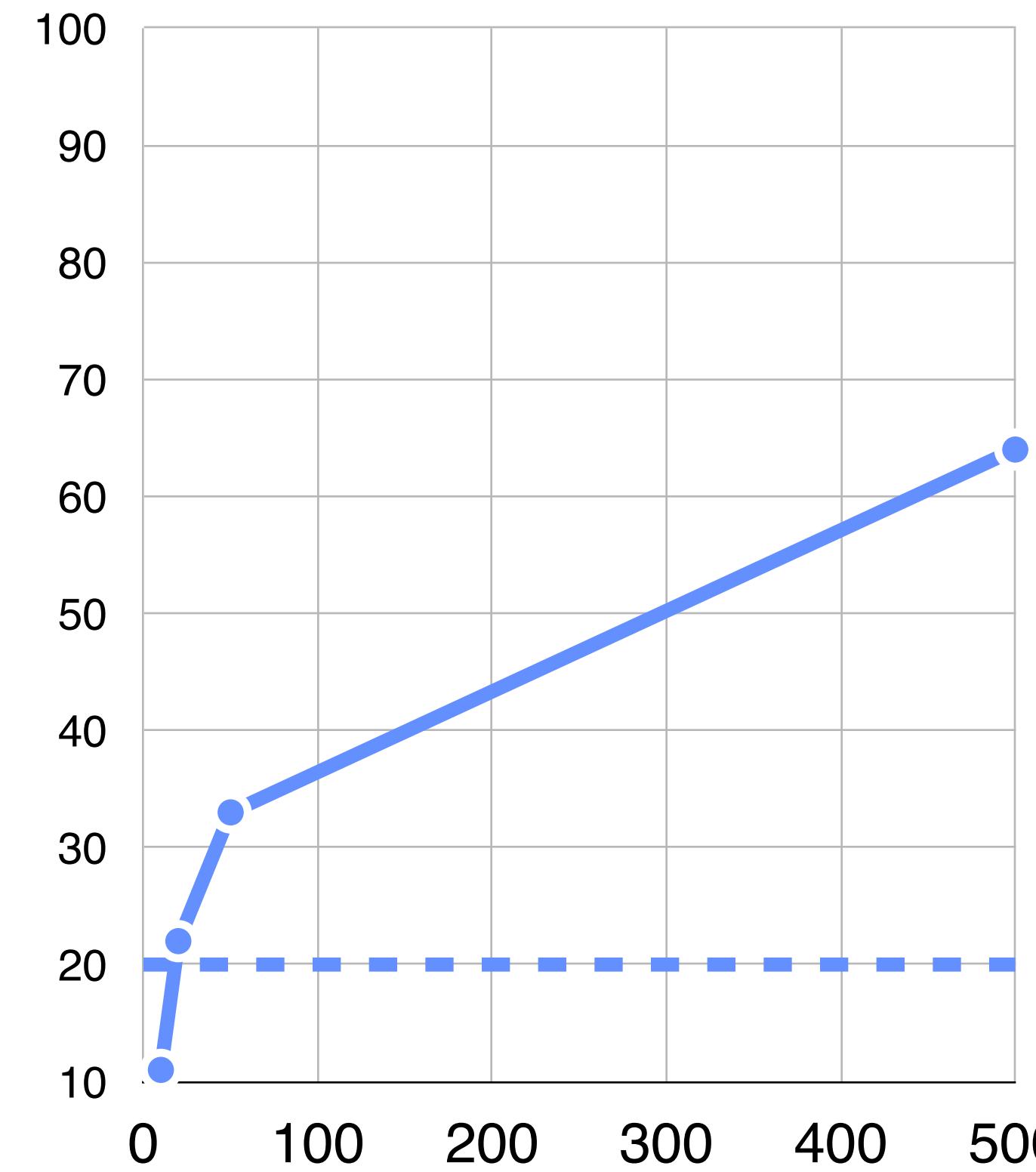
Text Classification

SST2 AG News Banking77



Entity Recognition

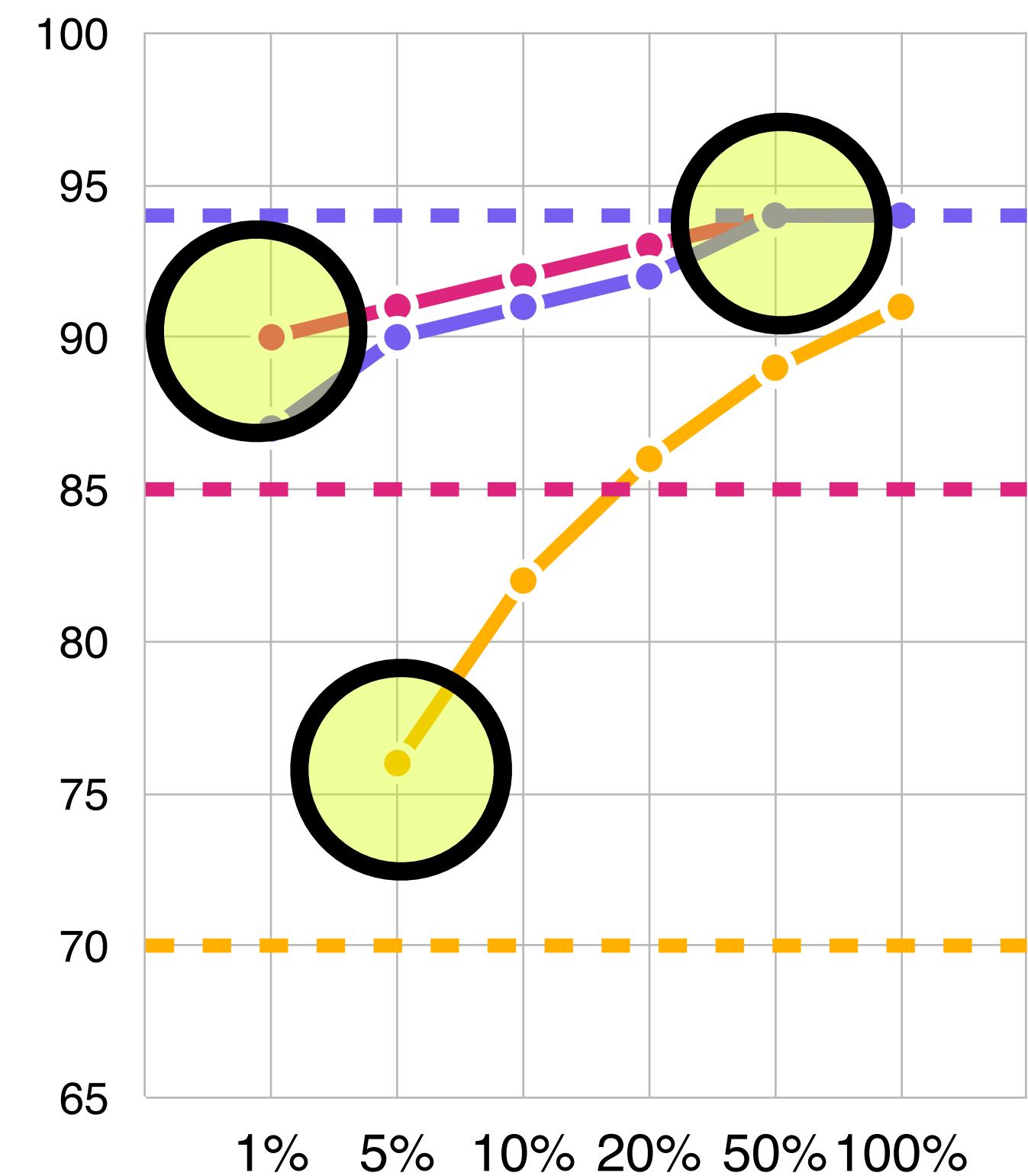
FabNER⁺⁺



EXPLORING
NOISE

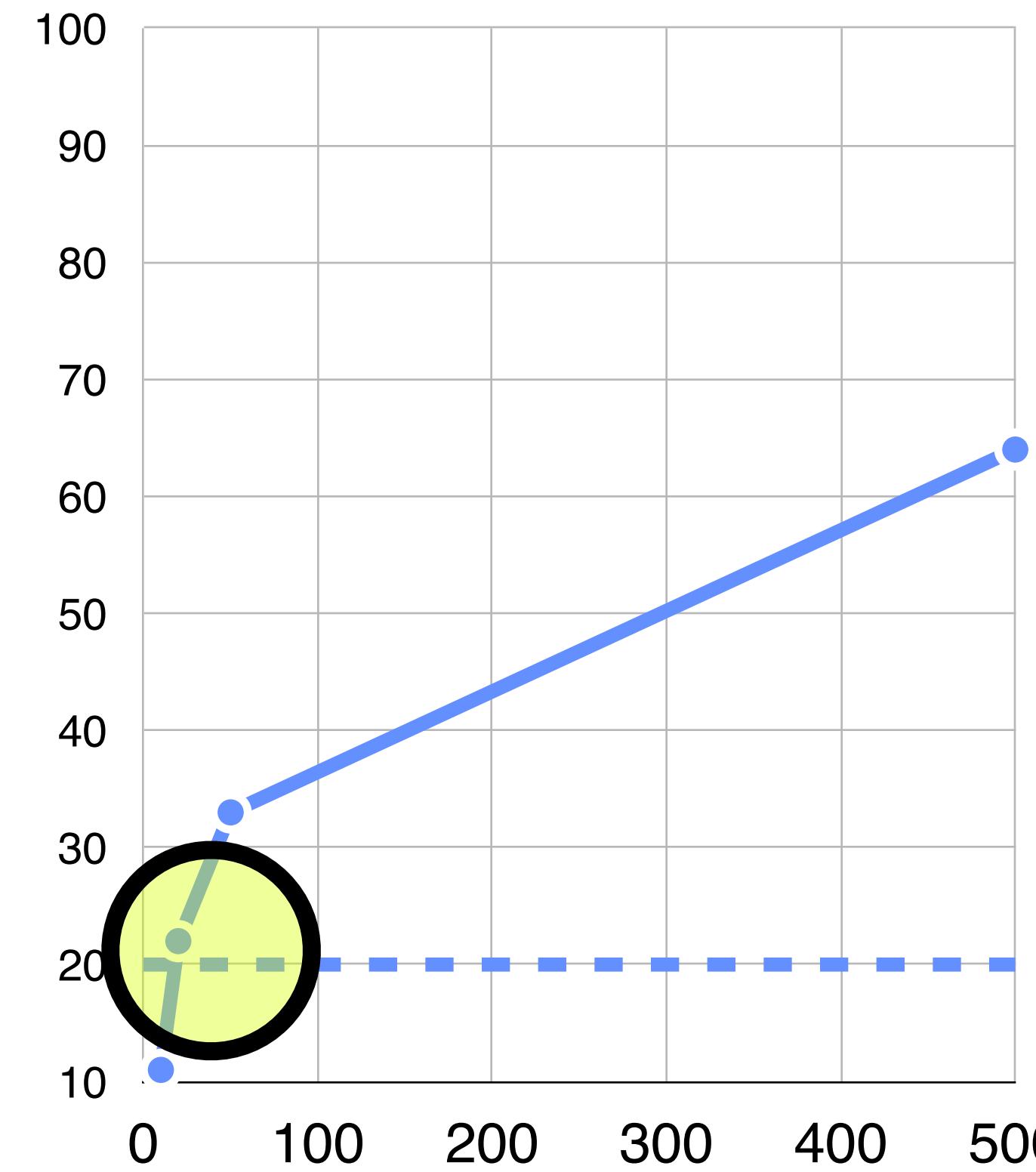
Text Classification

SST2 AG News Banking77



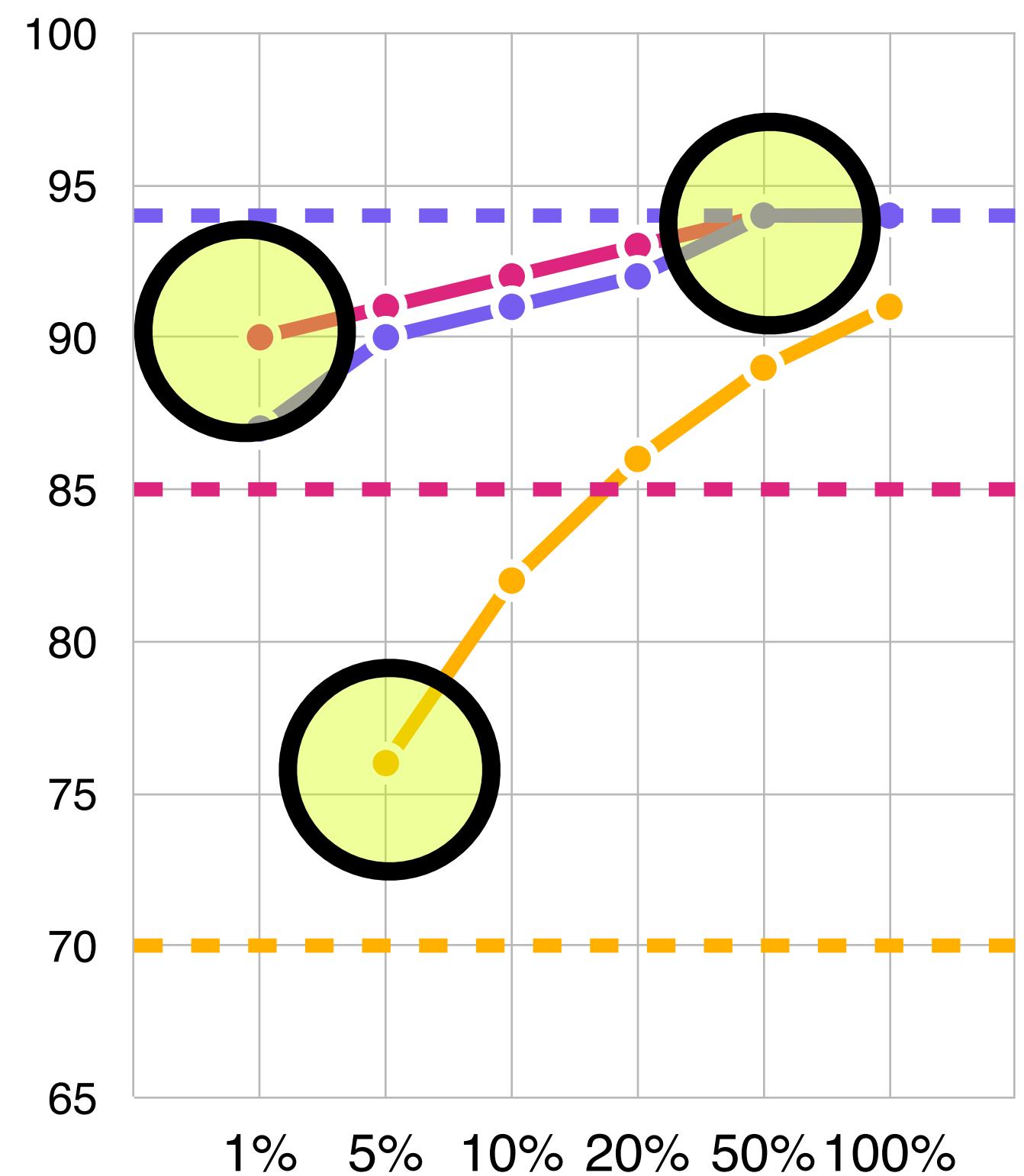
Entity Recognition

FabNER++



Text Classification

SST2 AG News Banking77



Entity Recognition

FabNER⁺⁺

