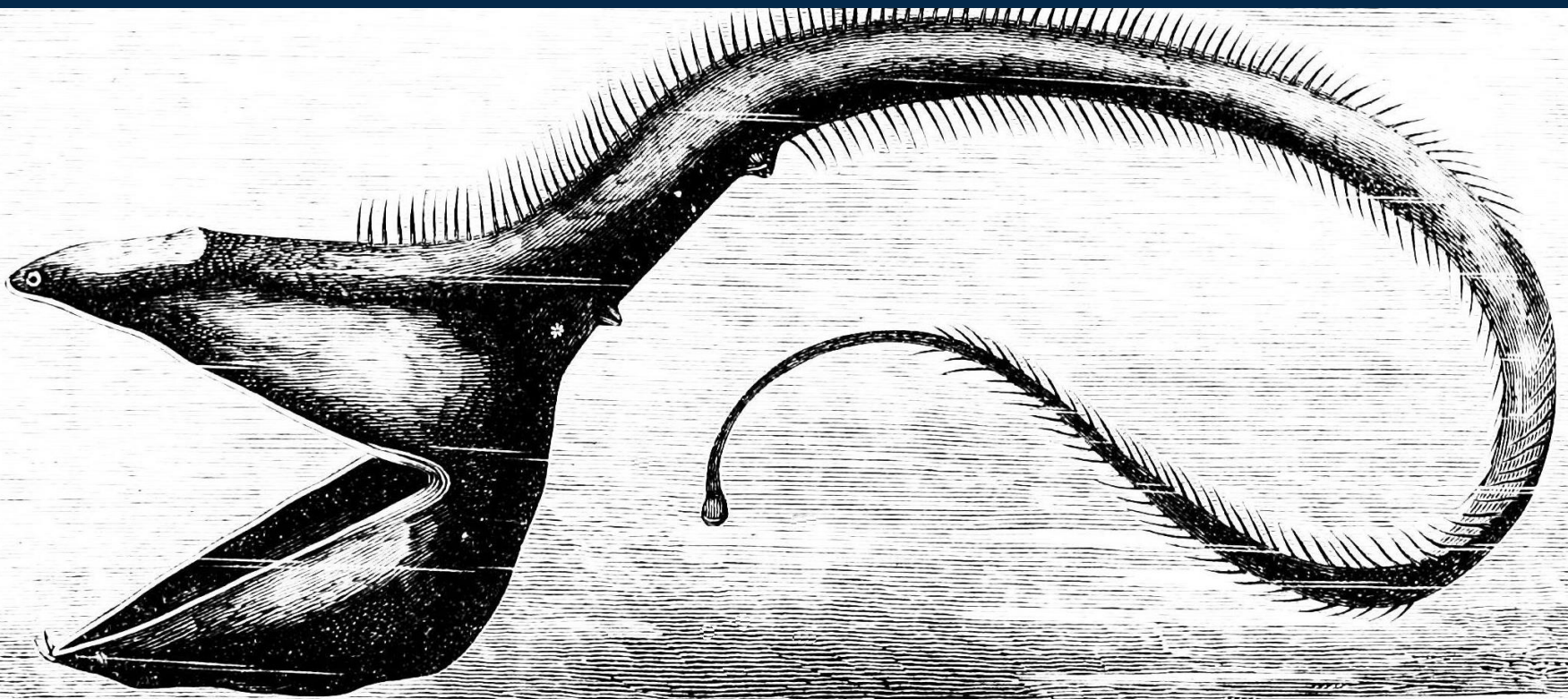


# Introducing *gobbli*

Deep learning with text doesn't have to be scary

Jason Nance  
Data Scientist  
RTI International



# Using a State of the Art Model for Text Classification

Line Count: 0 104 108 142 167 187

```
1 '''
2 Export our data in a format usable by Google's BERT model.
3 '''
4 import argparse
5 import dill
6 import pandas as pd
7 import numpy as np
8 from pathlib import Path
9 from collections import OrderedDict
10
11 from data import make_datasets
12
13 def preprocess_text(text_series):
14     return text_series
15
16 def
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

**BERT Evaluation**

Evaluate the performance of Google's BERT model.

```
In [1]: import pandas as pd
from sklearn.metrics import precision_recall_fscore_support, confusion_matrix

In [2]: input_df = pd.read_csv('/code/data/test_ids.csv')
predict_df = pd.read_csv('/code/data/bert/bert_output/test_results.tsv', sep='\t')
eval_df = input_df.merge(predict_df, left_index=True, right_index=True,
                           header=None, names=['p_false', 'p_true'])
eval_df.head()
```

id	p_true	p_false	precision	recall
0	932647	False	0.999875	0.000125
1	23245376	False	0.999875	0.000125
2	20500515	False	0.999875	0.000125
3	23565506	False	0.999875	0.000125
4	17510391	False	0.999875	0.000125

```
In [3]: def eval_threshold(threshold):
    pred_y = eval_df['p_true'] > threshold
    precision, recall, fscore, support = precision_recall_fscore_support(eval_df['y'], pred_y)
    print(f"Positive precision: {precision[1]}")
    print(f"Positive recall: {recall[1]}")
    print(f"Positive predictions: {pred_y.sum()} out of {pred_y.shape[0]}")
    print(f"-----")

for threshold in [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.75, 0.9, 0.95, 0.975, 0.98, 0.99]:
    eval_threshold(threshold)
```

```
Out [3]: For threshold 0.05:
Positive precision: 0.83338391582276176
Positive recall: 0.86162464989994398
Positive predictions: 659 out of 12655

For threshold 0.1:
Positive precision: 0.83462321792268928
```

```
1 train_df = pd.DataFrame(train_df_data)
2 del train_df_data
3 df_to_train_tsv(train_df, data_dir_path / "train.tsv")
4 del train_df
5
6 dev_df_data = []
7 for ex in valid_dataset:
8     dev_df_data.append({
9         "text": example_to_text(ex),
10        "label": ex.Relevant,
11    })
12
13 dev_df = pd.DataFrame(dev_df_data)
14 del dev_df_data
15
16 data_dir_path.mkdir(parents=True, exist_ok=True)
17
18 #!/bin/bash
19 set -u
20 echo "This script runs predictions from a trained BERT model."
21 echo "run 'train_bert.sh' first."
22
23 source bert_env.sh
24 export MODEL_CHECKPOINT=$BERT_BASE_DIR/bert_model.ckpt
25 export TRAINED_CLASSIFIER=$BERT_BASE_DIR/bert_classifier.py
26
27 docker run --rm --name=bert \
28     -v $(pwd)/data:/data \
29     -v $(pwd)/data/bert:/bert \
30     tensorflow/tensorflow:1.11.0-gpu-py3 python run_classifier.py \
31     --data_dir=/data \
32     --vocab_file=/bert/vocab.txt \
33     --bert_config_file=/bert/bert_config.json \
34     --init_checkpoint=/bert/bert_model.ckpt \
35     --max_seq_length=128 \
36     --train_batch_size=32 \
37     --learning_rate=5 \
38     --num_train_epochs=3.0 \
39     --output_dir=/bert/output_dir
```

**Data Export Script**

(Don't d)

```
1 #!/bin/bash
2 set -u
3 echo "This script will..."
4 echo "It expects export_b..."
5
6 if [[ -d data/bert ]]; then
7     git clone https://github.com/goog
8
9     cd ./data/bert
10    wget https://storage.googleapis.com/bert_model
11    unzip uncased_L-12_H-768_A-12.zip
12
13 fi
```

**Prediction Shell Script**

```
1 #!/bin/
2 set -u
3
4 echo "This script will..."
5 echo "It expects export_b..."
6
7 if [[ -d data/bert ]]; then
8     git clone https://github.com/goog
9
10    cd ./data/bert
11    wget https://storage.googleapis.com/bert_model
12    unzip uncased_L-12_H-768_A-12.zip
13
14 fi
```

**Training Shell Script**

```
1 #!/bin/
2 set -u
3
4 echo "This script will..."
5 echo "It expects export_b..."
6
7 if [[ -d data/bert ]]; then
8     git clone https://github.com/goog
9
10    cd ./data/bert
11    wget https://storage.googleapis.com/bert_model
12    unzip uncased_L-12_H-768_A-12.zip
13
14 fi
```














# Deep Learning: State of the Art

## Best Scores on DBpedia Classification Benchmark\*

- 2015: Char-level CNN (<https://arxiv.org/abs/1509.01626v3>)
- 2016: fastText (<https://arxiv.org/abs/1607.01759v3>), CNN/LSTM (<https://arxiv.org/abs/1602.02373v2>)
- 2017: DPCNN (<https://www.aclweb.org/anthology/papers/P/P17/P17-1052/>), M-ACNN (<https://arxiv.org/abs/1709.08294v3>)
- 2018: ULMFiT (<https://arxiv.org/abs/1801.06146v5>)
- 2019: BERT (<https://arxiv.org/abs/1810.04805>), MT-DNN (<https://arxiv.org/abs/1901.11504>), XLNet (<https://arxiv.org/abs/1906.08237>)

\* <https://paperswithcode.com/sota/text-classification-on-dbpedia>

# The Problem: Hard-Coding for Benchmark Problems

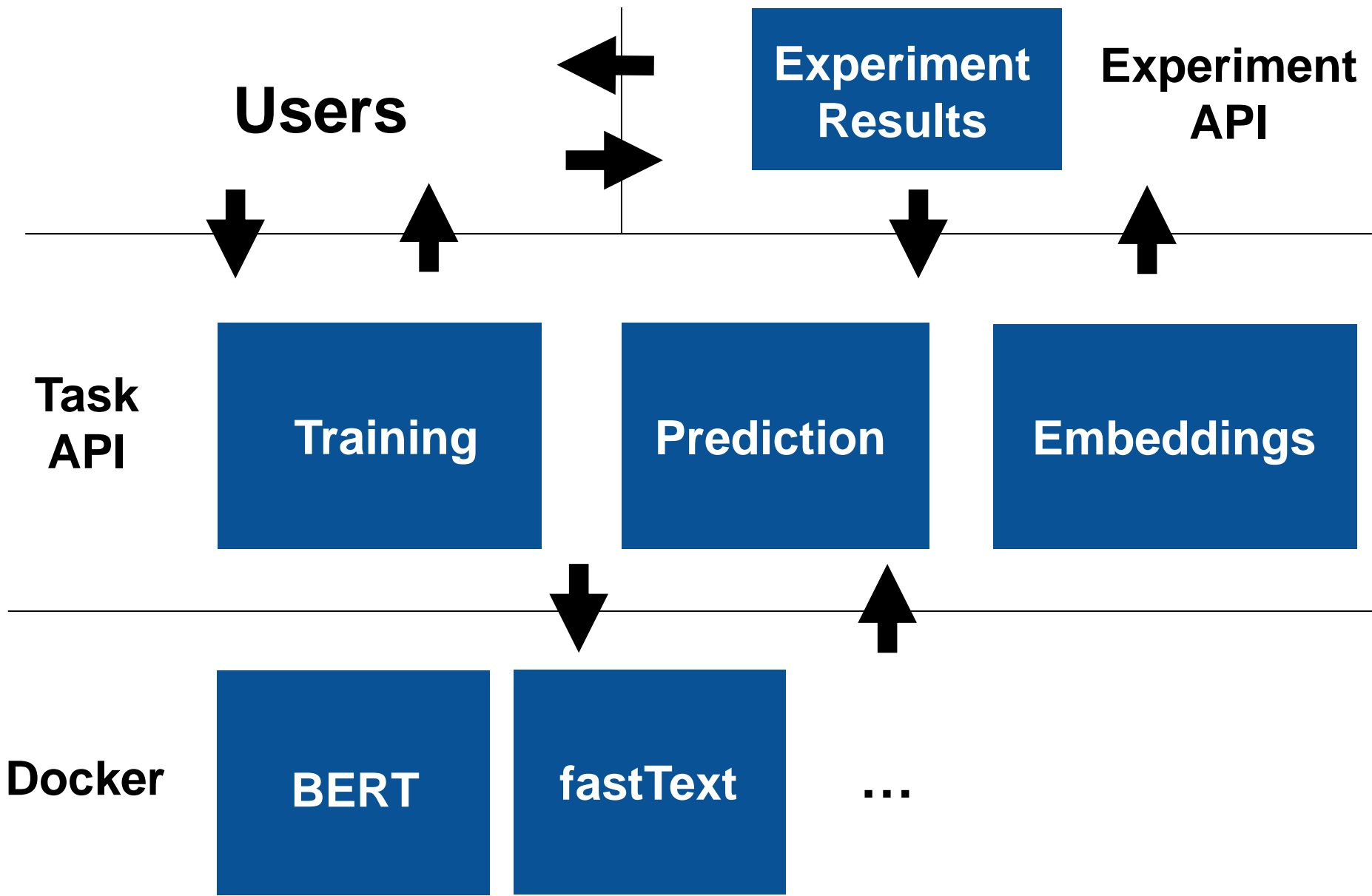
Rank	Name	Model	URL	Score	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m	MNLI-mm	QNLI	RTE	WNLI
1	Microsoft D365 AI & UMD	Adv-RoBERTa (ensemble)		88.8	68.0	96.8	93.1/90.8	92.4/92.2	74.8/90.3	91.1	90.7	98.8	88.7	89.0
2	Facebook AI	RoBERTa		88.5	67.8	96.7	92.3/89.8	92.2/91.9	74.3/90.2	90.8	90.2	98.9	88.2	89.0
3	XLNet Team	XLNet-Large (ensemble)		88.4	67.8	96.8	93.0/90.7	91.6/91.1	74.2/90.3	90.2	89.8	98.6	86.3	90.4
+	4	Microsoft D365 AI & MSR AI MT-DNN-ensemble		87.6	68.4	96.5	92.7/90.3	91.1/90.7	73.7/89.9	87.9	87.4	96.0	86.3	89.0
5	GLUE Human Baselines	GLUE Human Baselines		87.1	66.4	97.8	86.3/80.8	92.7/92.6	59.5/80.4	92.0	92.8	91.2	93.6	95.0
+	6	王玮		87.0	69.2	95.2	92.6/90.2	91.1/90.6	74.4/90.7	88.2	87.9	95.7	83.5	87.0
7	Stanford Hazy Research	Snorkel MeTaL		83.2	63.8	96.2	91.5/88.5	90.1/89.7	73.1/89.9	87.6	87.2	93.9	80.9	65.0
8	XLM Systems	XLM (English only)		83.1	62.9	95.6	90.7/87.1	88.8/88.2	73.2/89.8	89.1	88.5	94.0	76.0	71.0
9	Zhuosheng Zhang	SemBERT		82.9	62.3	94.6	91.2/88.3	87.8/86.7	72.8/89.8	87.6	86.3	94.6	84.5	65.0
10	Danqi Chen	SpanBERT (single-task training)		82.8	64.3	94.8	90.9/87.9	89.9/89.1	71.9/89.5	88.1	87.7	94.3	79.0	65.0
11	Kevin Clark	BERT + BAM		82.3	61.5	95.2	91.3/88.3	88.6/87.9	72.5/89.7	86.6	85.8	93.1	80.4	65.0
12	Nitish Shirish Keskar	Span-Extractive BERT on STILTs		82.3	63.2	94.5	90.6/87.6	89.4/89.2	72.2/89.4	86.5	85.8	92.5	79.8	65.0
13	Jason Phang	BERT on STILTs		82.0	62.1	94.3	90.2/86.6	88.7/88.3	71.9/89.4	86.4	85.6	92.7	80.1	65.0
14	廖亿	RGLM-base (Huawei Noah's Ark Lab)		81.0	55.1	94.2	90.7/87.7	89.5/88.7	72.2/89.4	85.6	85.1	92.1	78.5	65.0
+	15	Jacob Devlin		80.5	60.5	94.9	89.3/85.4	87.6/86.5	72.1/89.3	86.7	85.9	92.7	70.1	65.0

<https://gluebenchmark.com/leaderboard/>

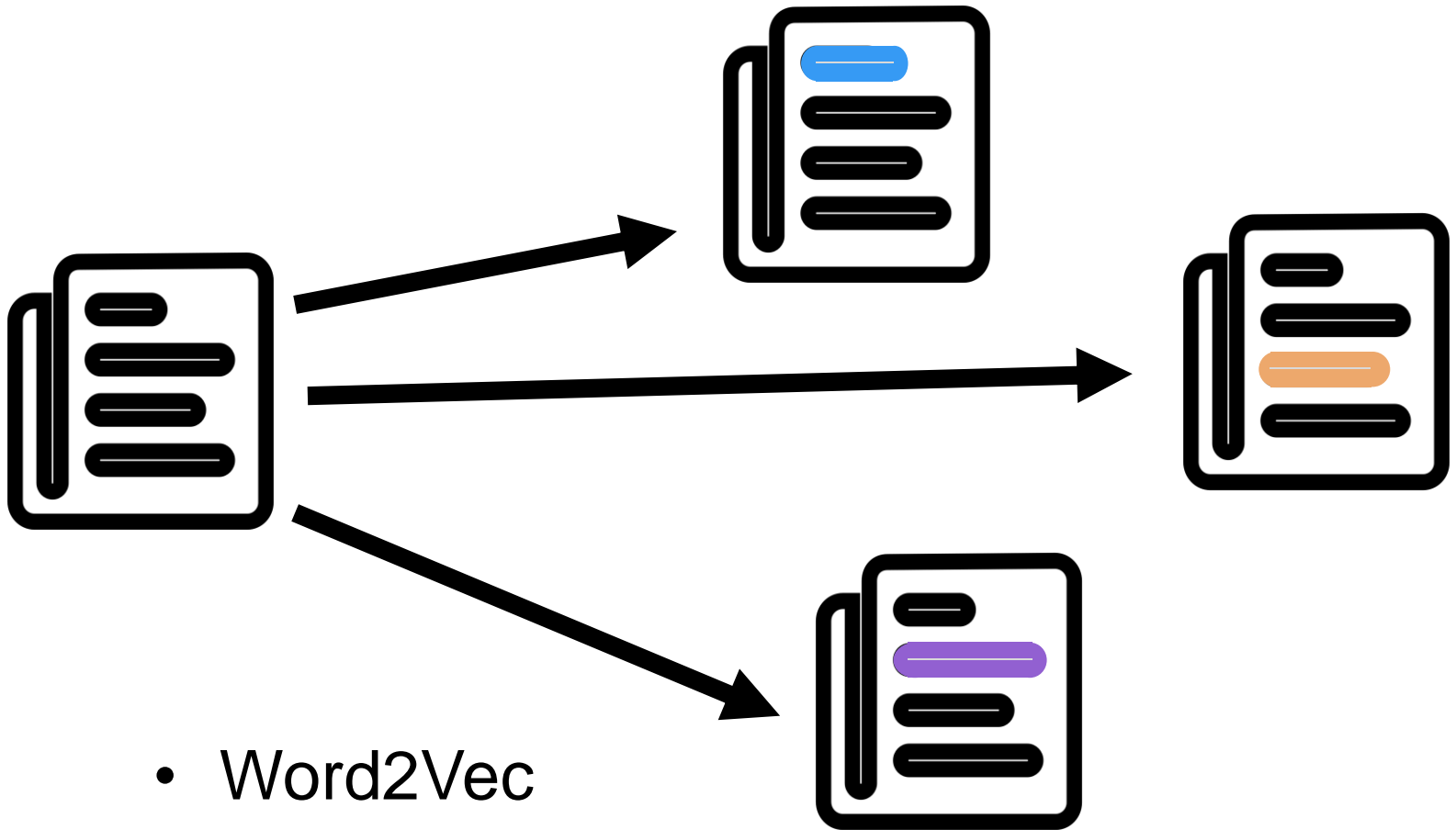
# gobbli: A Uniform Interface for Text Deep Learning Models



# gobbli: Library Design



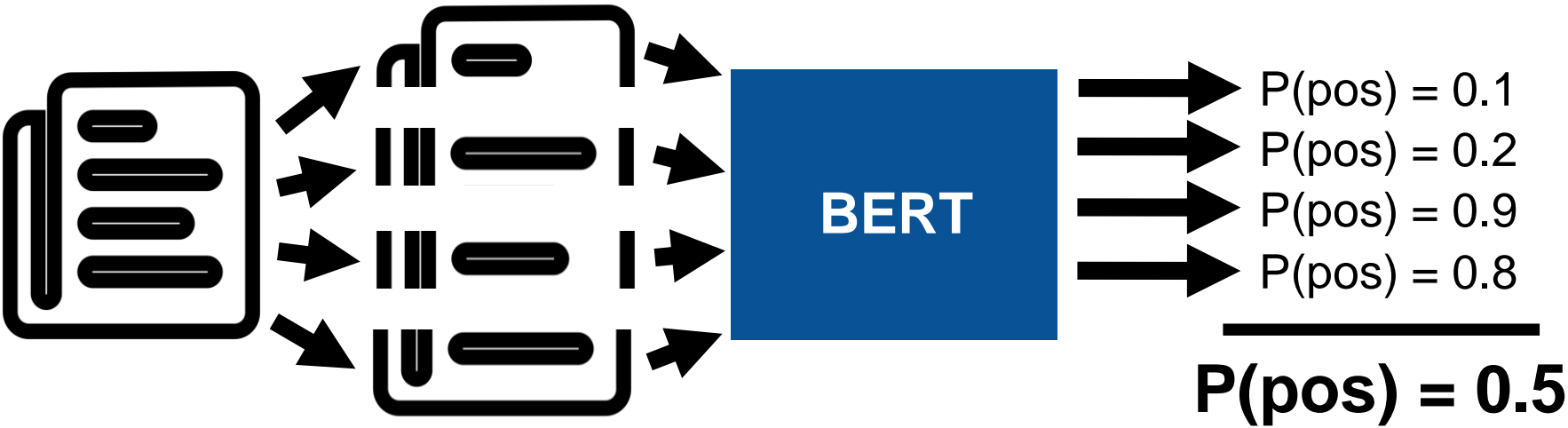
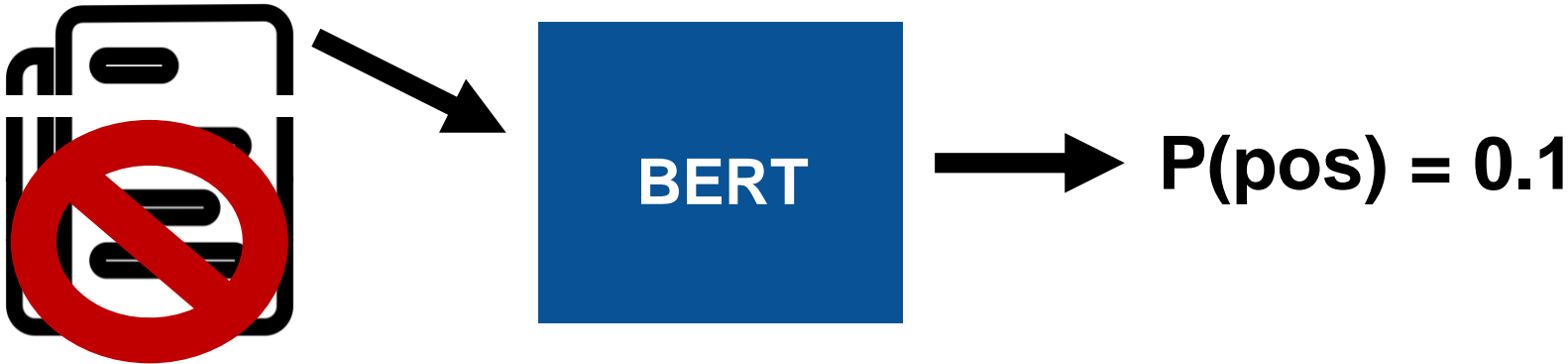
## Data Augmentation



- Word2Vec
- WordNet
- BERT Masked Language Model

# gobbli: Additional Features

## Document Windowing







- + Cross-platform
- + Abstracts dependency management
- Latency/overhead



- + Parallel/distributed training
- Experiment API only

# Example Experiment Results

# Example Experiment Results: Metrics

Metrics:

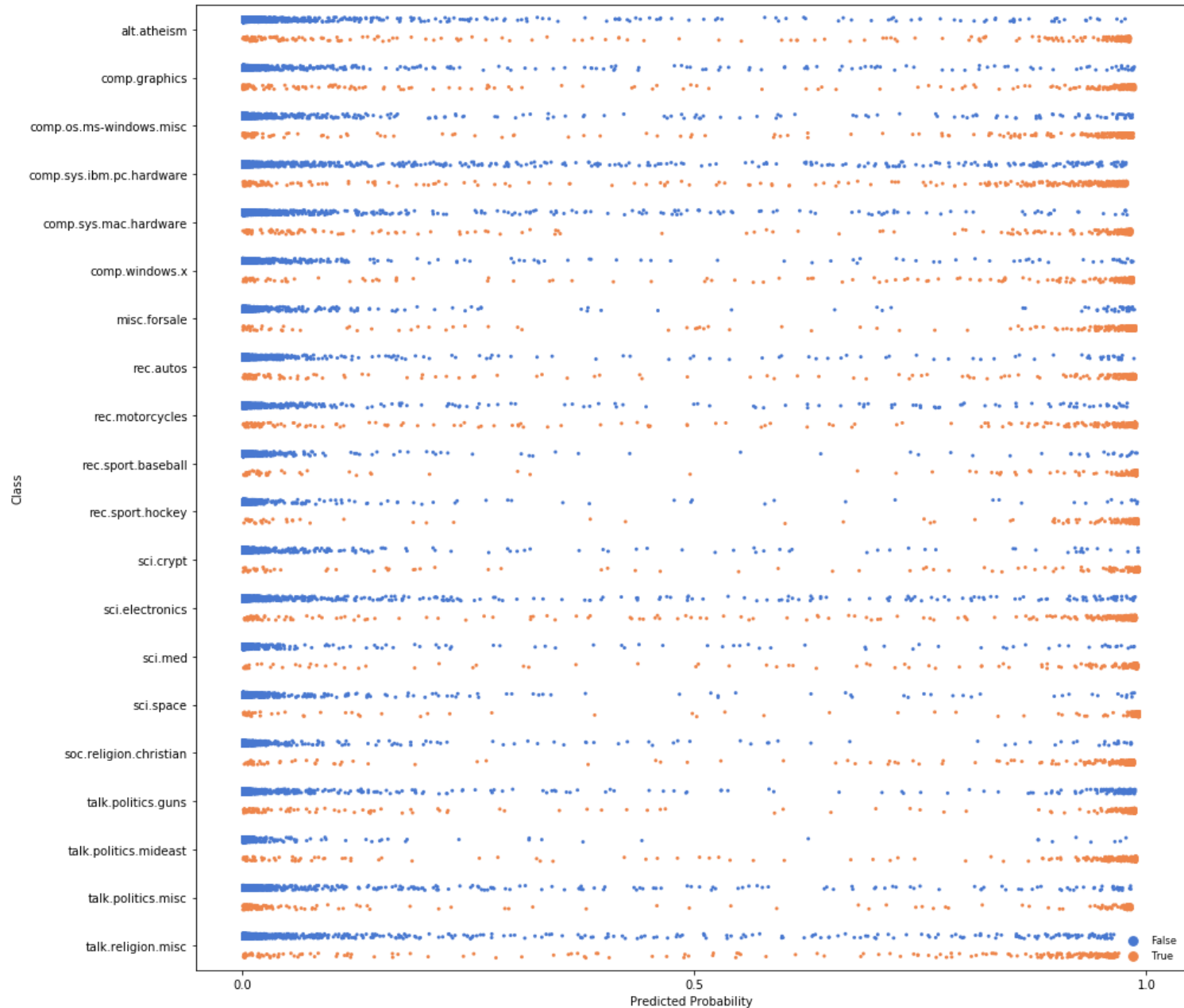
-----  
Weighted F1 Score: 0.8806791429898766  
Weighted Precision Score: 0.8806909370983464  
Weighted Recall Score: 0.88068  
Accuracy: 0.88068

Classification Report:

-----

	precision	recall	f1-score	support
neg	0.88	0.88	0.88	12500
pos	0.88	0.88	0.88	12500
accuracy			0.88	25000
macro avg	0.88	0.88	0.88	25000
weighted avg	0.88	0.88	0.88	25000

# Example Experiment Results: Plot



# Example Experiment Results: Errors Report

**True Class:** comp.os.ms-windows.misc

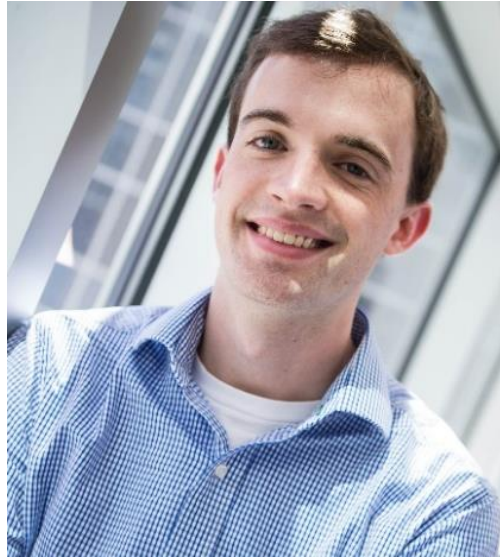
**Predicted Class:** sci.med (Probability: 0.97)

**Text:**

“My wife is a physiotherapist and she is looking for some cliparts of skeleton and male/female body. We're currently using Windows Draw which can import all kind of graphic formats. Therefore, anything will do. Please advise ...”

# gobbli: Status and Next Steps

- Initial open source release on GitHub
  - <https://github.com/RTIInternational/gobbli/>
  - Models implemented: BERT, MT-DNN, USE, fastText, pytorch\_transformers (XLNet, XLM, BERT, RoBERTa)
- Next steps:
  - Support multilabel classification
  - Helper module for downstream tasks using embeddings
  - Helper module for exploratory descriptives
  - Other bug fixes/enhancements requested by the community



**Jason Nance**  
**Data Scientist**  
**RTI International**  
jnance@rti.org