

Learning Robust Representations using a Change Point Framework

Anonymous Authors

ABSTRACT

This work proposes change point embeddings (CAPE), word embeddings learned using a novel change-point objective framework. A unique feature of this objective function is that it trains representations to be optimal at distinguishing between contexts. We gain improved separation between context representations, and therefore achieve improved performance in the presence of domain shift - a change in the data distribution between an algorithm's training dataset, and a dataset it encounters when deployed. Without needing fine-tuning, this change point framework outperforms state-of-the-art models in clustering topics, and is also more robust against noise in topic classification. Code and trained models are available (anonymized) [at this URL](#).

CCS CONCEPTS

• **Computing methodologies** → **Natural language processing**; **Information extraction**.

KEYWORDS

Word embeddings, Bayesian change-point model, robustness

ACM Reference Format:

Anonymous Authors. 2021. Learning Robust Representations using a Change Point Framework. In *IGKDD Conference on Knowledge Discovery and Data Mining (KDD '21)*, June 00–01, 2021, Virtual Event. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Word embedding — the mapping of words into numerical vector spaces — has become a basic building block in many natural language processing (NLP) tasks. Many embedding techniques have been developed, including pretrained embedding techniques (e.g. Word2Vec [26], GloVe [28] and fastText [25]). More recently, deep neural network based models (e.g. ELMO [29], BERT [11]) have emerged. These models are trained to encode complex word relationships, leading to improved performance in many NLP tasks.

The training of these architectures often takes place at the word and sentence level. Word masking is a commonly-used training strategy, where models are trained by running a prediction task that either predicts a word after seeing its surrounding context of words or predicts a context after seeing one member word [11, 24]. The resulting representation is particularly good at tasks at the

word and sentence levels, such as mapping a context down to a single word [34]. However, despite the addition of objectives aimed at understanding relationships between sentences [11], topic-level information is still underutilized in the training of word representations.

We propose a novel embedding training method, change point embedding (CAPE), based on a probabilistic change-point model that optimizes the embedding's ability to distinguish between different topics. The key idea of this method is to construct training passages with a known transition of topics and then model the change of topics using a change point model. By training the embeddings to optimize the detection of the change point, it incorporates topic information into the learning process, gearing the embedding towards learning topic-discriminative features.

The use of this statistical model results in optimally separated word representations. As well-separated representations are less likely to be corrupted by noise, it improves the tolerance of the representation against corruption introduced by noise. This is especially useful for language models, as the majority of the world's languages do not have reliable textual or expert resources [14] and language itself changes between different domains and over time.

In systematic comparisons, we show that CAPE substantially improves performance in unsupervised tasks such as identification of topic transitions, including situations with domain shift, and improves robustness to noise in text classification. We also assess the performance on NLP tasks that are not topic-related using the GLUE benchmark, evaluating the types of transfer learning that are possible with CAPE's framework. In addition, we also develop novel strategies for evaluating pre-trained language models in topic discovery, differentiation and robustness to noise.

2 RELATED WORK

2.1 Text representation

Word embeddings associate each word with a vector. Due to their ability to capture syntactic and semantic information in words, pre-trained word vectors are a core component in current NLP architectures. Numerous embeddings are available, for example Word2Vec [24, 26], GloVe [28] and fastText [25]. Word2Vec and FastText are trained using CBOW, a word prediction strategy, while GloVe is trained on a combination of global word-word co-occurrence statistics and local context window methods. One drawback of these techniques is that the embedding representation is context independent, i.e. the vector representing a word stays fixed regardless of surrounding words. This makes them less adaptive to changes in the language environment.

More recently, deep neural network encoders have been shown to outperform embeddings across a broad range of diverse NLP tasks [30, 34]. These models, such as BERT [11], RoBERTa [23], XLNet [38], GPT-3 [7], demonstrate the efficacy of transformer models [33], in combination with millions to billions of parameters,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).

KDD 2021, June 00–01, 2021, Virtual Event

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

extensive computational resources, and large amounts of training data [17].

By training them using word prediction strategies [20][15] in a large general corpus, these techniques generate contextual representations of input tokens that are infused with information of its neighborhood. The same word can have very different representations in different situations, thus taking context into account. After training, the full neural network models that were used to train these embeddings can be used for downstream tasks, unlike pre-trained embeddings which require additional structures (e.g. classifiers) to make predictions. Relatively simple fine-tuning can be performed on these general models with task-specific data to improve performance in downstream tasks [18].

2.2 Change point analysis

Change point analysis is a powerful tool for detecting whether any changes take place in a data stream, and to locate where the changes occur [4]. Bayesian Online Changepoint Detection algorithm (BOCD) [3] identifies change points by modeling the probability distribution of the “run length”, i.e. the elapsed time since the most recent change point. By leveraging conjugate-exponential models, it achieves efficient, exact and online Bayesian inference. It has been applied and extended in numerous ways, for example to include variational inference for non-exponential distributions [32], time lags that use future information to better detect change points [8], and non-stationary spatio-temporal processes [21].

3 THE MODEL

In the CAPE encoder model, each word w_i is represented using θ_{w_i} where θ_{w_i} is a d -dimensional vector with an associated covariance matrix $\Sigma_{w_i} = \text{var}(\theta_{w_i})$. Now suppose a sentence j consists of words w_1, \dots, w_n . We encode it using (x_j, Σ_j) as follows:

$$\begin{aligned} x_j &= \frac{1}{n}(\theta_{w_1} + \dots + \theta_{w_n}) \\ \Sigma_j &= \text{var}(x_j) = \frac{1}{n^2}(\Sigma_{w_1} + \dots + \Sigma_{w_n}) \end{aligned} \quad (1)$$

where

$$\Sigma_{w_i} = \begin{bmatrix} \sigma_{w_i,1}^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_{w_i,d}^2 \end{bmatrix}_{d \times d}.$$

Here we assume all θ_{w_i} are independent, and all Σ_{w_i} diagonal to enforce independence across different dimensions of θ_{w_i} .

We model x_j using a multivariate normal distribution [27] and use a normal prior over the unknown parameter μ as follows:

$$\begin{aligned} x_j | \mu, \Sigma_j &\sim \mathcal{N}(\mu, \Sigma_j) \\ \mu &\sim \mathcal{N}(\mu_0, \Sigma_0) \end{aligned} \quad (2)$$

To learn representations that incorporate topic-level information, we develop a training strategy that consists of passages with known transitions of topics, and formulate the learning process into a model-based change-point framework. We then construct and optimize an objective function, based on BOCD, to learn $(\hat{\theta}, \hat{\Sigma})$ that both recognize and differentiate topics.

Below in 3.1 we give a brief overview of BOCD. We then define the change-point based objective function in 3.2 and show that it renders less correlated (3.2.2) representations.

3.1 An overview of the BOCD algorithm

The BOCD algorithm [3] is used to find ‘change points’, i.e. locations of abrupt changes in a sequence of observations, x_1, \dots, x_T , in an online manner. It assumes that the time series can be partitioned into a sequence of non-overlapping segments, where the boundaries between segments are the change points. Given a segment, observations within the segment ρ are assumed to be conditionally independently and identically distributed from a segment-specific probability distribution $P(x_t; \eta_\rho)$.

Briefly, BOCD estimates change points by probabilistically estimating the “run length” distribution r_t , i.e. the estimated number of observations passed at time t since the last change point. If a change point occurs at t , $r_t = 0$; otherwise, $r_t = r_{t-1} + 1$, i.e. the run length grows. Specifically, denote the observations that have been seen so far as $x_{1:t} = \{x_1 \dots x_t\}$, then the run length distribution at the t^{th} observation is $R_{t,r} = P(r_t = r | x_{1:t}) \propto P(r_t, x_{1:t})$, where $r \in [0, 1, \dots, t-1]$.

The joint distribution $P(r_t, x_{1:t})$ is modeled in terms of $P(r_{t-1}, x_{1:t-1})$ as follows,

$$P(r_t, x_{1:t}) = \begin{cases} \sum_{r_{t-1}} \frac{1}{\lambda} P(x_t | r_{t-1}, x_{1:t-1}; \eta_t^{(r_{t-1})}) P(r_{t-1}, x_{1:t-1}) & \text{if } r_t = 0. \\ (1 - \frac{1}{\lambda}) P(x_t | r_{t-1}, x_{1:t-1}; \eta_t^{(r_{t-1})}) P(r_{t-1}, x_{1:t-1}) & \text{if } r_t = r_{t-1} + 1. \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

where $\frac{1}{\lambda}$ is the prior probability of have a change point at time t and the posterior predictive distribution $P(x_t | r_{t-1}, x_{1:t-1}; \eta_t^{(r_{t-1})})$ is specified according to the data. Given $R_{t,r}$, one can estimate the most probable run length at t as $\text{argmax}_r R_{t,r}$.

To estimate $R_{t,r}$, the algorithm processes the data stream in an online fashion. For each newly observed data point x_t , the algorithm uses the previously computed $P(r_{t-1} | x_{1:t-1})$ and the current data point x_t to estimate $P(r_t | x_{1:t})$ for all the possible run lengths ($r_t \in [0, 1, \dots, t-1]$). The algorithm then updates all the parameters to take into account x_t , in the standard Bayesian manner, and the posterior is used as a prior for the next observed data point. When $P(x_t)$ is chosen to be a member of exponential family and conjugate priors are used, the parameters can be updated efficiently by conjugacy. Afterwards the algorithm returns the estimated run length probabilities, $R = \{R_{t,r} : t \in [0, 1, \dots, T], r_t \in [0, 1, \dots, t-1]\}$ and the estimated most probable change points. As $t^2/2 R_{r,t}$ ’s need to be computed, the computational complexity of BOCD is $O(t^2)$.

3.2 The change-point based objective function and learning framework

To learn representations that incorporate topic-level context information, we construct a training set where each observation contains two different contexts, and optimize an objective function based on BOCD that tries to differentiate the contexts.

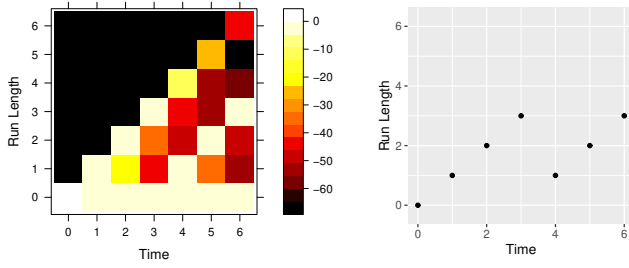


Figure 1: Visualization of the run length distributions for a simulated example of 6 observations ($m = 3$), with a change point at the 4th point. Left: Estimated log run-length probabilities, R . The color of the square at (t, r) corresponds to $\log P(r_t = r | x_{1:t})$. Right: Maximum a posteriori (MAP) run length estimates, $\max_r R_{t,r}$, which suggests the data can be divided into two segments: $\{x_1, x_2, x_3\}$ and $\{x_4, x_5, x_6\}$.

For the training set, we randomly choose two different contexts from a very large corpus (e.g. two different Wikipedia paragraphs) and then concatenate a string of m observations from each context (e.g. 2 sentences from each paragraph). We then encode them using (θ, Σ) . This creates a data stream with a known change point at the $(m+1)$ th observation. Using BOCD, we want to realize posterior run-length distributions that reflect the correct location of the change point, so we update (θ, Σ) to improve performance at the objective function below.

$$obj = R_{1,1} + R_{2,2} + \dots + R_{m,m} + R_{m+1,1} + R_{m+2,2} + \dots + R_{2m,m} \quad (4)$$

Figure 1 illustrates the case of $m = 3$. A context change occurs between the first three and the last three observations. Thus $obj = R_{1,1} + R_{2,2} + R_{3,3} + R_{4,1} + R_{5,2} + R_{6,3}$. The true run length increases over the first three observations, resets to 0, and increases again over the last three, so we optimize the probability estimates to match. For example, $R_{6,3}$ is maximized because the most probable run-length estimate at $t = 6$ should be $r_6 = 3$.

Large values of (4) are achieved when the first m observations are represented as a compact cluster and the last m observations map to a distant compact cluster. Thus, maximizing (4) with respect to (θ, Σ) will force the embedding to learn representations that best differentiate contexts. The full process is summarized in Algorithm 1. As all terms in (4) can be obtained from one pass of BOCD, the computational complexity for (4) is the same as BOCD, i.e. $O(m^2)$ for one passage. Other choices of $m(m \geq 2)$ are possible and a larger m would lead to a longer runtime.

Note that this task does not require human annotation; the paragraph information is already in Wikipedia and in most text datasets. This is known as self-supervision, where we use a supervised loss, but the data provides its own supervision. This allows us to use the vast amounts of unlabeled text that are available.

Algorithm 1 Training algorithm

```

1: procedure TRAIN
2:   procedure INITIALIZE
3:     Randomly initialize the encoder parameters  $(\theta, \Sigma)$ 
4:   for K iterations do
5:     Pick 2 distinct contexts in the training set and  $m$  obser-
       vations from each context.
6:     Encode the  $2m$  observations using  $\theta, \Sigma$ . ▷ Eq 1
7:     Get  $R$  using BOCD and calculate the objective. ▷ Eq 3, 4
8:     Improve  $obj$  (wrt  $\theta, \Sigma$ ) using stochastic gradient meth-
       ods.
9:   return trained encoder parameters  $(\hat{\theta}, \hat{\Sigma})$ 

```

3.2.1 Additional model specification and hyperparameter optimization. By standard Bayesian conjugate computation in (2), the posterior predictive distribution $P(x_t | r_{t-1}, \mathbf{x}_t^{(r_{t-1})}, \eta_t^{(r_{t-1})})$ in (3) is a multivariate Normal distribution with run-specific parameters $\eta_t^{(r)} = \{\mu_t^{(r)}, \Sigma_t^{(r)}\}$. Below we give the initial parameters $\eta^{(0)} = \{\mu^{(0)}, \Sigma^{(0)}\}$ and Bayesian update procedures [27]. As in [3], we start with $R_{00} = P(r_0 = 0) = 1$.

$$\begin{aligned} \mu^{(0)} &= (0, 0 \dots 0)_d \\ \Sigma^{(0)} &= \begin{bmatrix} (\sigma^2)^{(0)} & & \\ & \ddots & \\ & & (\sigma^2)^{(0)} \end{bmatrix}_{d \times d} \\ \Sigma_{t+1}^{(r+1)} &= ((\Sigma_t^{(r)})^{-1} + \Sigma_t^{-1})^{-1} \end{aligned} \quad (5)$$

$$\mu_{t+1}^{(r+1)} = \Sigma_{t+1}^{(r+1)} (\Sigma_t^{-1} x_t + (\Sigma_t^{(r)})^{-1} \mu_t^{(r)})$$

In this model the hyperparameters are

- The embedding size, d
- The number of sentences being compared, m
- The optimization algorithm (gradient descent variety), and its learning rate
- The batch size for gradient descent
- The strategy used to initialize (θ, Σ) . We used a random-uniform initialization, where the ranges are hyperparameters.
- $(\sigma^2)^{(0)}$, the prior variance within BOCD.
- The constant hazard parameter λ within BOCD.

In this paper we present experimental results for $d \in \{10, 100, 500\}$. The number of sentences compared at each iteration was fixed at $m = 2$. Experiments with $m = 3$ found no appreciable difference except the models took longer to converge. The batch size was set at 25 for all experiments. Experiments showed that AdaGrad [13] was the best optimization algorithm. The parameter λ was fixed at 10. Preliminary experiments show that the choice of λ too made little difference.

For each d , the other hyperparameters were tuned using random search [6]. We picked random sets of hyperparameters, ran the algorithm for $5e4$ iterations with each set, then picked the hyperparameter set with the highest objective value at the end.

Other architectural choices include the preprocessing and tokenization strategies in Section 5.1. We used fairly straightforward tokenization techniques.

3.2.2 Separation of embedding vectors representing different contexts . We outline a brief argument to show that the optimization of (4) will separate vectors representing different contexts.

As described earlier in the setup of (4), where $m = 3$, the true run length increases over the first three observations, then resets to 0 and increases again over the last three observations. This means $R_{3,3} = P(r_3 = 3|x_{1:3})$ and $R_{4,1} = P(r_4 = 1|x_{1:4})$ will be maximized. Hence, $R_{4,4} = P(r_4 = 4|x_{1:4})$ will be minimized, since $\sum_{r=1}^4 R_{4,r} = 1$. Consequently, $P(r_4, x_{1:4})$ will be small and $P(r_3, x_{1:3})$ will be large. By Eq 3 (middle case) for $t = 4$, $P(x_4|r_3 = 3, \mathbf{x}^{(3)} = x_{1:3}; \eta_4^{(3)})$ must be very small.

Based on the conjugate specification in (2) and (5), the posterior predictive distribution $P(x_4|r_3 = 3, \mathbf{x}^{(3)} = x_{1:3}; \eta_4^{(3)})$ is a multivariate Normal distribution with parameters found through conjugate Bayesian updating [27] over $x_{1:3}$ and $\Sigma_{1:3}$. It takes a small value when its kernel $(x_4 - \mu_4^{(3)})^T (\Sigma_4^{(3)})^{-1} (x_4 - \mu_4^{(3)})$ is large. Note that,

$$\begin{aligned} (x_4 - \mu_4^{(3)})^T (\Sigma_4^{(3)})^{-1} (x_4 - \mu_4^{(3)}) & \quad \text{Mahalanobis distance} \\ = (x_4 - \mu_4^{(3)})^T P \Lambda^{-1} P^T (x_4 - \mu_4^{(3)}) & \quad \text{eigendecomposition} \quad (6) \\ = \mathbf{b}^T \Lambda^{-1} \mathbf{b} & \quad \mathbf{b} = P^T (x_4 - \mu_4^{(3)}) \quad (7) \\ = \sum_i \frac{b_i^2}{\lambda_i} & \quad (8) \end{aligned}$$

So as $R_{4,4}$ is minimized through gradient descent, the kernel above is maximized, and it is preferential to move x_4 away from $\mu_4^{(3)}$ primarily in the direction associated with the smallest eigenvalue λ_i of $\Sigma_4^{(3)}$, because that yields the greatest increase in statistical (Mahalanobis) distance. This eigenvector is orthogonal to all the other eigenvectors, which together capture the major sources of variation in the distribution fit on $x_{1:3}$. Thus, the estimated x_4 is nearly orthogonal with $x_{1:3}$. As the algorithm converges, the representations of dissimilar contexts are distant and decorrelated.

Similar reasoning shows that representations in similar contexts (i.e. each of $\{x_1, x_2, x_3\}$, and $\{x_4, x_5, x_6\}$) will be close in Mahalanobis (and Euclidean) distance, and therefore more correlated. Together, this indicates that the optimization of (4) will separate vectors representing different contexts.

4 EVALUATION OF EMBEDDINGS

We assess the embeddings in both topic-related tasks and other NLP tasks. For topic-related tasks, we evaluate the coherence of topics (4.1) and robustness to noise in topic classification (4.2). For other NLP tasks (4.3), we use the General Language Understanding Evaluation (GLUE) benchmark ([34]). In addition, we also evaluate distances between sentence vectors (4.4), which implies discriminability between sentences. All evaluations are performed on the test datasets, encoded using trained embeddings.

4.1 Clustering assessment

Representations that successfully differentiate between different contexts are expected to have low separation between representations within the same contexts, and high separation between dissimilar contexts. We quantify the separation quality using modularity.

Modularity, originally developed for community detection, is designed to measure how well a graph separates into a specific partitioning [10]. Here we create a passage with $m = 3$ sentences each from 2 unique paragraphs. We represent each sentence in a data stream as a node in the graph, and define the weight of the edge between node u and node v as $A_{uv} = \text{corr}(x_u, x_v)$, where x_u and x_v are the respective encoded representations and $\text{corr}(x_u, x_v)$ is computed using Pearson correlation. We then calculate modularity, given by

$$Q = \frac{1}{|E|} \sum_{uv} \left[A_{uv} - \frac{k_u k_v}{|E|} \right] \delta(c_u, c_v), \quad (9)$$

where A_{uv} is the weight of the edge between nodes u and v , k_u and k_v are the sum of edge weights associated with nodes u and v respectively, $\delta(c_u, c_v)$ equals 1 iff u and v originate from same context, and zero otherwise, and $|E|$ is the total sum of edge weights in the graph, i.e. the sum of all values in the upper triangle of the adjacency matrix A . This measure computes the normalized difference between the observed edge weights A_{uv} and the expected weights by chance $\frac{k_u k_v}{|E|}$ for all pairs of observations from the same context reflecting the strength of association beyond chance. Q is in the range $(-1, 1)$. A higher value means that the embeddings demonstrate (relatively) strong connections between sentences in the same topic and weak connections between sentences in different topics.

We found modularity to work well in practice, reflecting separation between topics. Appendix A demonstrates how modularity reflects the separation between topics in an example passage.

4.2 Robustness evaluation in classification tasks

The robustness of word embeddings against noise in the data is critical for generalizability and reproducibility. We evaluate the robustness of the trained word embeddings in the task of topic classification when the text to be classified consists of added noise.

To proceed, we split a labeled dataset into training and test sets, and add a fixed amount of random noise (random words) to each element of the test set. After encoding the training and test sets using the trained word embeddings, we then train a classifier using the training set, and run the classifier on the test set and report the classifier’s accuracy as a function of the amount of added noise. (Note: when evaluating BERT, we used its internal classification capabilities.) This experiment measures how robust a classifier trained on the representation of clean data is against noisy data. When no noise is added, this task reduces to the commonly reported accuracy metrics.

4.3 GLUE benchmark

The General Language Understanding Evaluation (GLUE) benchmark ([34]) is a collection of diverse natural language understanding tasks, each with an associated dataset. We use these reference datasets, commonly encountered in NLP, to understand what kinds of transfer learning are possible with this training framework, and how CAPE trades other strengths with topic discriminability. A detailed description of the tasks is in the Appendix (Table 5).

We encode the dataset associated with each task using trained word embeddings, train a classifier on the designated training data, then make predictions on the designated test dataset. The results are scored by the GLUE evaluation server.

4.4 Dispersion of sentence vectors

Dispersion of sentence vectors in a dataset reflects how well sentence vectors span the embedding space. A higher value implies that sentence vectors are well separated, offering a higher tolerance to corruption introduced by noise.

To quantify vector dispersion in a dataset, we first define a normalized distance statistic. Given a collection of sentences and a sentence vector v from a dataset, we define the normalized distance from v to the collection as

$$\text{Normalized Distance}(v) = \frac{\text{median distance from } v \text{ to other sentence vectors in the dataset}}{\sqrt{\sum_i v_i^2}} \quad (10)$$

To measure the overall vector dispersion in a dataset, we sample random vectors and sentence collections in the dataset, then compute the normalized distance for each random vector, and report the median normalized distance. Specifically in this paper we repeat the following process for 100 times: 1) we sample a group of 50 paragraphs in a dataset, 2) pick out one sentence as v , 3) calculate the median of pairwise distances from v to all other sentences in the 50 paragraphs, and 4) compute the normalized distance value using (10). The median of the 100 normalized distance values serves as an overall measure of vector dispersion in a dataset.

5 EMPIRICAL STUDY SETUP

5.1 Datasets

We downloaded Wikipedia (July 2019) and cleaned it using a modified version of Matt Mahoney’s pre-processing script¹. This left us with 310M total paragraphs. Of these, we kept the subset that have at least 3 sentences with 4+ words each, obtaining 24M paragraphs. The last 300K paragraphs were split off to make the Wikipedia test set, and the rest were used to make the training data. This training set had a full vocabulary of size 4.8M, but to save time and space we restricted training and testing to the 400K most common words in the training set. Each paragraph is treated as a distinct context.

To test the embedding’s performance under domain shift, we selected Bookscorpus as another test set. Bookscorpus consists of large datasets of free books on the internet written by yet unpublished authors [19, 22, 39, 40]. Its narrative style is quite different from Wikipedia training data. For testing, we randomly picked out

one book from this collection and cleaned it similar to Wikipedia, resulting in 1904 paragraphs.

For the robustness analysis, we picked the AG News corpus [1]. It consists of news articles collected from various sources online, each annotated by news type. For simplicity, we kept only articles classified into 3 large classes: ‘Sci/Tech’, ‘Sports’, and ‘World’. Following [16], we concatenated each article’s title and short description together (average length 42 words), and used that to predict the news type. We randomly sampled 3200 articles for the robustness analysis, splitting off 3000 articles for training the classifier and 200 for testing.

5.2 Training process

We maximized the objective function on the Wikipedia training set using Tensorflow’s [2] implementation of an adaptive gradient descent algorithm (AdaGrad [13]). We trained 3 models corresponding to $d \in \{10, 100, 500\}$, denoted as CAPE(d), each for 7-8M iterations. The procedure took {39, 40, 60} hours respectively on an i7 8700K CPU.

5.3 Testing process

We evaluate CAPE against three competing methods: the fast-Text 300-dimensional English embeddings (2018) [25], the 300-dimensional Wikipedia+Gigaword GloVe vectors [28], and the 1024 dimensional BERT embeddings (BERT-Large, Uncased, Whole Word Masking), extracted using bert-as-service [37]. We run the following analyses:

- (A) We randomly sample 1000 paragraphs from Wikipedia and Bookscorpus respective test sets and compute normalized distance (4.4) to evaluate vector dispersion in these two datasets. For each dataset, we also construct passages with 2 distinct contexts ($m=2$) and calculate modularity (4.1) to evaluate topic discrimination (avg runtime: < 5m). For Bookscorpus, we first perform the analysis without fine tuning, then we fine-tune the CAPE encoder on a 100-paragraph subset of Bookscorpus for improving performance. In the fine-tuning, we set up initial encoder parameters starting at the previously trained values ($\hat{\theta}$, $\hat{\Sigma}$) and all other hyperparameters unchanged as in the original training. In each iteration we present 2 distinct contexts and $m = 2$ observations from each context. Fine tuning was done until convergence of the objective, usually for less than 1e5 iterations (avg runtime: 35m), compared to 8e6 in the pretraining process in (5.2).
- (B) We run the classification test (4.2) on the AG news dataset, evaluating resistance to added noise. We train a SVM classifier on the clean data, and add a fixed amount (0-100) of random noise words at random positions to each article in the test set. We then assess the classifier’s accuracy for classifying the article’s news type as a function of the amount of noise (avg runtime: < 10m per test).
- (C) For the GLUE benchmark, we train a classifier to the encoded training data, then evaluated results on the GLUE server. We use random forest because it scales better than SVM despite slightly lower accuracy. This was done without fine-tuning the encoder.

¹<http://mattmahoney.net/dc/textdata>

6 RESULTS AND DISCUSSION

6.1 CAPE generates better topic discrimination and higher vector separation in the test sets

We evaluated the modularity and normalized distance on the test sets of Wikipedia passages and Bookcorpus passages (without fine tuning). Table 1 shows a numerical summary of the results. Figures 4-5 in Appendix B show the full distributions.

CAPE has significantly higher median modularities than the other methods, suggesting that it clusters paragraphs (i.e. topics) much more coherently. It also has significantly higher median normalized distance. This suggests that it has better separated sentence embedding, indicating a higher efficiency to span the embedding space. It is also observed that its modularities and normalized distance increase with the increase of its dimension. Though Bookscopus has very different narrative style from the training set (Wikipedia), CAPE’s superiority in performance maintains, even without fine-tuning. In contrast, BERT does not seem to transfer well to these tasks.

	Modularity		Normalized distance	
	Wikipedia	Bookscopus	Wikipedia	Bookscopus
CAPE (10)	0.304	0.163	1.173	0.910
CAPE (100)	0.396	0.286	1.300	1.243
CAPE (500)	0.465	0.376	1.389	1.343
Glove	0.128	0.107	0.714	0.652
FastText	0.146	0.124	0.820	0.763
BERT	0.097	0.094	0.537	0.534

Table 1: Summary of clustering evaluation results in Wikipedia and Bookscopus. The medians of modularity and normalized distance are reported. No fine tuning was applied to Bookscopus.

To evaluate the effectiveness of our fine tuning procedure, we fine-tuned the encoder parameters on the Bookscopus dataset. From the set of 1904 paragraphs, we split off 100 paragraphs for fine-tuning and the rest for modularity and dispersion evaluation. We performed additional training of the normal CAPE training algorithm on the 100 paragraphs for 1e5 iterations. Table 2 shows a significant improvement in unsupervised performance, achieved in 1% of the original pre-training time.

6.2 CAPE is robust to noise in text classification

Next, we evaluated CAPE in text classification using the AG news data. We focused on the robustness of the embeddings against the distortion detailed in Section 4.2. Note that we only trained the classifier and did not fine-tune the encoder. When training BERT, we ran a grid search over the fine-tuning options recommended in the BERT paper (Appendix 3), varying batch size, learning rate and number of epochs.

	Modularity		Normalized distance		Elapsed Time
	Before	After	Before	After	
CAPE (10)	0.163	0.33	0.910	1.40	28m
CAPE (100)	0.286	0.42	1.243	1.41	30m
CAPE (500)	0.376	0.44	1.343	1.41	44m

Table 2: CAPE fine-tuning results on Bookscopus. The median normalized distances and modularities on Bookscopus before and after fine-tuning are displayed.

As shown in Fig 2 (top), though all the three word embedding methods have broadly similar performance when no noise is added, the CAPE text classification has a much smaller accuracy reduction when noise is added. BERT’s results (Fig 2 bottom) vary across fine-tuning settings. The overall relationship between their classification accuracy and the amount of added noise largely follow the general pattern as for FastText and GloVe, showing a larger accuracy reduction with added noise than CAPE. This indicates that CAPE is much more robust against this noise model than all the other three methods. With higher model dimension, CAPE shows improved classification accuracy in both the cases with and without added noise. CAPE’s robustness is likely due to its property of generating more separated sentence vectors (6.1). Larger distance between sentence vectors offers higher tolerance to added noise.

6.3 GLUE benchmark results

We evaluated performance on the GLUE benchmark for a variety of NLP tasks. Table 3 shows the results as scored by the evaluation server². BERT outperforms CAPE in all these tasks. This is not surprising, as BERT has many more parameters, and is trained on much more data and with much more computational resources than CAPE. The performance of a language model is known to scale as a power-law of model size, dataset size, and the amount of computation [7, 17].

However, we do find evidence of transfer learning from CAPE’s original training task to these tasks, especially those involving semantics (e.g. QQP, MRPC and RTE). The most promising CAPE performance is in the Microsoft Research Paraphrase Corpus (MRPC) [12]. In this task, the aim is to identify if two sentences are paraphrases of each other, and this might be viewed as similar to distinguishing contexts. On the other end of the scale, CAPE does badly at CoLA [35] because it is good at ignoring noise. In this task, each example is a sequence of words annotated with whether it is a grammatical English sentence or not.

7 CONCLUSION

In this paper we propose a novel change-point based framework for training and evaluating embeddings. CAPE, with its distributional assumptions and optimized separation through a change-point objective, strongly outperforms conventional models such as fastText,

²<https://gluebenchmark.com/leaderboard>

	MNLI (m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE
BERT	86.7/85.9	89.3	92.7	94.9	60.5	86.5	89.3	70.1
CAPE (100)	48.6/48.9	70.8	67.0	70.9	7.3	56.1	79.9	54.2
CAPE (500)	48.8/49.4	72.3	69.1	71.9	10.4	57.3	79.5	55.0

Table 3: GLUE Test results, scored by the evaluation server. F1 scores are reported for MRPC, Spearman correlations are reported for STS-B, Matthew’s correlation is reported for CoLA, and accuracy scores are reported for the other tasks.

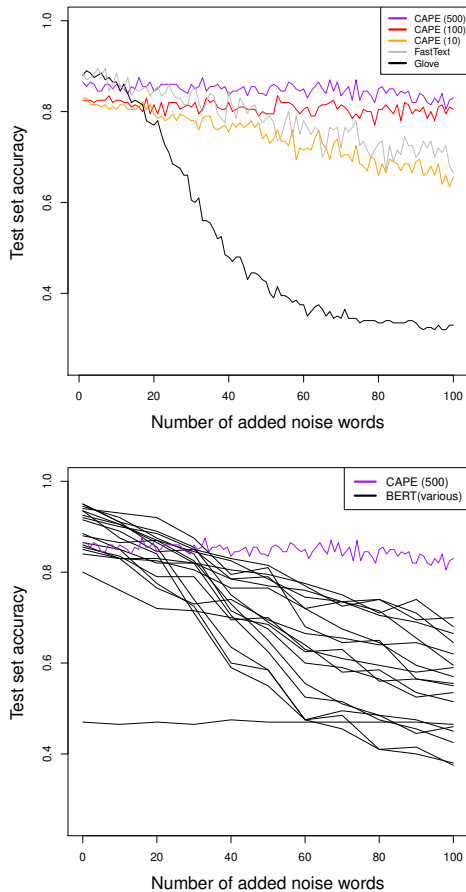


Figure 2: Accuracy in text classification as a function of added noise in the AG news dataset. Each line shows how a classifier trained (or fine-tuned) on a clean training set performs when evaluated on a test set with added noise. Top: CAPE, FastText and GloVe. Bottom: BERT with various fine-tuning settings.

GLoVe or BERT at distinguishing between topics in the unsupervised setting.

The CAPE framework also demonstrates robust performance in topic classification when noise is present. We attribute CAPE’s robustness to the improved separation between the word embeddings,

as seen using the normalized distance statistic. Other methods have closely packed vectors that are more likely to overlap, so they are more sensitive to distortions in the texts. Robustness is essential for generality and reproducibility of word embeddings in NLP tasks, as such tasks often involve changing language environments and adaptability to different contexts is highly desirable.

CAPE does not match BERT’s accuracy in the tasks in the GLUE benchmark, but the use of a larger more diverse pre-training dataset and more embedding dimensions is expected to improve CAPE accuracy.

REFERENCES

- [1] [n.d.]. AG News corpus. http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html. Accessed: 2020-02-06.
- [2] Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. 2016. Tensorflow: a system for large-scale machine learning. In *OSDI*, Vol. 16. 265–283.
- [3] Ryan Prescott Adams and David JC MacKay. 2007. Bayesian online changepoint detection. *arXiv preprint arXiv:0710.3742* (2007).
- [4] Samaneh Aminikhanghahi and Diane J Cook. 2017. A survey of methods for time series change point detection. *Knowledge and information systems* 51, 2 (2017), 339–367.
- [5] Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. 2009. The Fifth PASCAL Recognizing Textual Entailment Challenge. (2009).
- [6] James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *The Journal of Machine Learning Research* 13, 1 (2012), 281–305.
- [7] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).
- [8] Michael Byrd, Linh Nghiem, and Jing Cao. 2017. Lagged Exact Bayesian Online Changepoint Detection. *arXiv preprint arXiv:1710.03276* (2017).
- [9] Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. *arXiv preprint arXiv:1708.00055* (2017).
- [10] Mingming Chen, Konstantin Kuzmin, and Boleslaw K Szymanski. 2014. Community detection via maximization of modularity and its variants. *IEEE Transactions on Computational Social Systems* 1, 1 (2014), 46–65.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [12] William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the International Workshop on Paraphrasing*.
- [13] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12, Jul (2011), 2121–2159.
- [14] Ewan Dunbar, Xuan Nga Cao, Juan Benjumea, Julien Karadayi, Mathieu Bernard, Laurent Besacier, Xavier Anguera, and Emmanuel Dupoux. 2017. The zero resource speech challenge 2017. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 323–330.
- [15] Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. Learning distributed representations of sentences from unlabelled data. *arXiv preprint arXiv:1602.03483* (2016).
- [16] Di Jin, Zhijiang Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. Is bert really robust? natural language attack on text classification and entailment. *arXiv*

- preprint arXiv:1907.11932 (2019).
- [17] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* (2020).
 - [18] Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882* (2014).
 - [19] Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2015. Skip-Thought Vectors. *arXiv preprint arXiv:1506.06726* (2015).
 - [20] Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. In *Advances in neural information processing systems*. 3294–3302.
 - [21] Jeremias Knoblauch and Theodoros Damoulas. 2018. Spatio-temporal Bayesian On-line Changepoint Detection with Model Selection. *arXiv preprint arXiv:1805.05383* (2018).
 - [22] Sosuke Kobayashi. 2018. Homemade BookCorpus. <https://github.com/soskek/bookcorpus>.
 - [23] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
 - [24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
 - [25] Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhres, and Armand Joulin. 2018. Advances in Pre-Training Distributed Word Representations. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*.
 - [26] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 746–751.
 - [27] Kevin P Murphy. 2007. Conjugate Bayesian analysis of the Gaussian distribution. *def 1, 2 σ^2* (2007), 16.
 - [28] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
 - [29] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365* (2018).
 - [30] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of EMNLP (Austin, Texas)*. Association for Computational Linguistics, 2383–2392.
 - [31] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of EMNLP*. 1631–1642.
 - [32] Ryan D Turner, Steven Bottone, and Clay J Stanek. 2013. Online variational approximations to non-exponential family change point models: with application to radar tracking. In *Advances in Neural Information Processing Systems*. 306–314.
 - [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762* (2017).
 - [34] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. In the *Proceedings of ICLR*.
 - [35] Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2018. Neural Network Acceptability Judgments. *arXiv preprint 1805.12471* (2018).
 - [36] Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of NAACL-HLT*.
 - [37] Han Xiao. 2018. bert-as-service. <https://github.com/hanxiao/bert-as-service>.
 - [38] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *arXiv preprint arXiv:1906.08237* (2019).
 - [39] Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning Books and Movies: Towards Story-like Visual Explanations by Watching Movies and Reading Books. In *arXiv preprint arXiv:1506.06724*.
 - [40] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning Books and Movies: Towards Story-Like Visual Explanations by Watching Movies and Reading Books. In *The IEEE International Conference on Computer Vision (ICCV)*.

A SEPARATION BETWEEN TOPICS IN AN EXAMPLE PASSAGE

Table 4: A sample passage of six sentences with two topics

Topic: Feral cats

A feral cat is a cat that lives outdoors and has had little or no human contact

They do not allow themselves to be handled or touched by humans, and will run away if they are able

They typically remain hidden from humans, although some feral cats become more comfortable with people who regularly feed them

Topic: Atlanta

Atlanta is the capital of, and the most populous city in, the US state of Georgia

With an estimated 2017 population of 486,290, it is also the 39th most populous city in the United States

The city serves as the cultural and economic center of the Atlanta metropolitan area, home to 5.8 million people and the ninth largest metropolitan area in the nation.

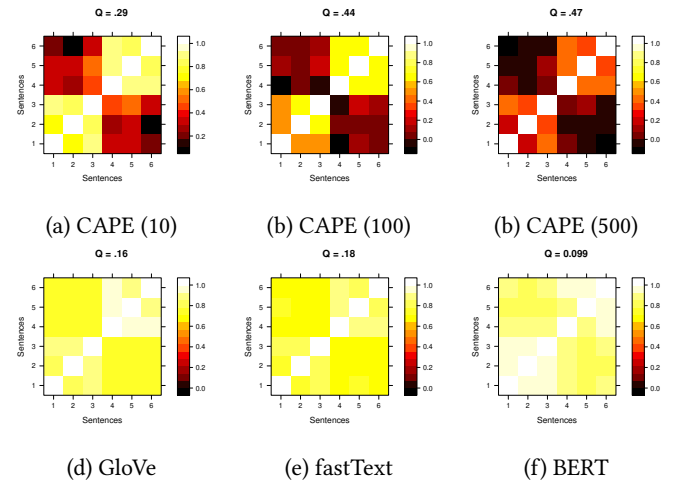


Figure 3: Heatmaps show the correlations between all pairs of sentences from Table 4.

As an illustration, we compare the models on a Wikipedia passage (Table 4). The first three sentences are about feral cats, and the next three are about the city of Atlanta. Figure 3 shows the results.

In this example, CAPE achieves a much higher modularity score ($Q = .29, .44, .47$) than the other three word embedding methods ($Q = 0.18, 0.16$, and 0.14 for fastText, GloVe, and BERT, respectively), suggesting that CAPE has much better topic coherence and topic separation. CAPE also shows improved performance with increasing dimension.

B BOXPLOTS OF MODULARITY AND NORMALIZED DISTANCE IN WIKIPEDIA AND BOOKSCORPUS TEST SET

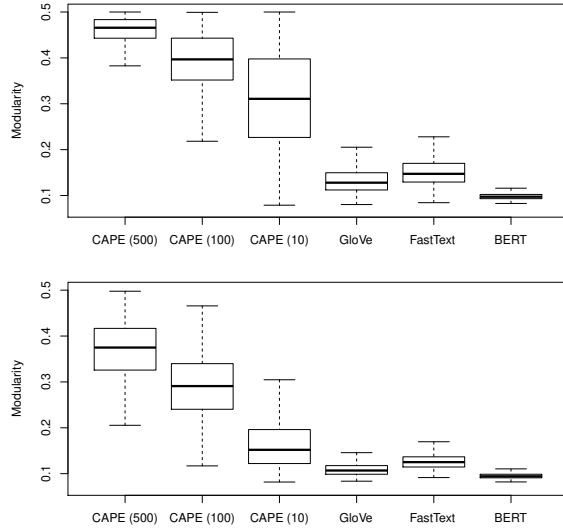


Figure 4: Distribution of modularity values in the Wikipedia (top) and Bookscorpus (bottom) test set.

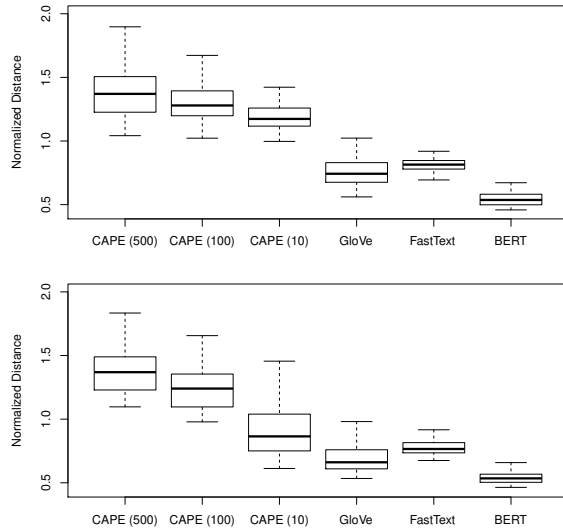


Figure 5: Distribution of normalized distance values in the Wikipedia (top) and Bookscorpus (bottom) test set.

Table 5: Descriptions and examples of the tasks in the GLUE benchmark.

Task	Label
MNLI [36]	Given a pair of sentences, predict whether the second sentence is an entailment, contradiction, or neutral with respect to the first one.
	"Conceptually cream skimming has two basic dimensions - product and geography."
	"Product and geography are what make cream skimming work. "
QQP ³	Determine if two questions asked on Quora are semantically equivalent.
	"How is the life of a math student? Could you describe your own experiences?"
	"Which level of prepration is enough for the exam jlt5?"
QNLI [30, 34]	Given a (question, sentence) pair, determine whether or not the answer to the question lies in the sentence.
	"When did the third Digimon series begin?"
	"Digimon Tamers takes a darker and more realistic approach to its story featuring Digimon who do not reincarnate after their deaths and more complex character development in the original Japanese."
SST-2 [31]	Given a single sentence, determine the sentiment.
	"saw how bad this movie was "
	"the greatest musicians "
CoLA [35]	Predict whether or not an English sentence is linguistically "acceptable".
	"We yelled ourselves."
	"Our friends won't buy this analysis, let alone the next one we propose."
STS-B [9]	This task contains sentence pairs annotated [1-5] according to how similar the sentences are.
	"A plane is taking off."
	"An air plane is taking off. "
MRPC ⁴ [12]	Given a sentence pair, determine whether the sentences in the pair are semantically equivalent [0,1].
	"Amrozi accused his brother , whom he called the witness , of deliberately distorting his evidence . "
	"Referring to him as only the witness , Amrozi accused his brother of deliberately distorting his evidence . "
RTE [5]	Similar to MNLI, but with much less training data.

³<https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs>

⁴<https://www.microsoft.com/en-us/download/details.aspx?id=52398>