# Context is Everything: Finding Meaning Statistically in Semantic Spaces

#### Eric Zelikman

Stanford University ezelikman@stanford.edu

# **Abstract**

This paper introduces Contextual Salience (CoSal), a simple and explicit measure of a word's importance in context which is a more theoretically natural, practically simpler, and more accurate replacement to tf-idf. CoSal supports very small contexts (20 or more sentences), out-of context words, and is easy to calculate. A word vector space generated with both bigram phrases and unigram tokens reveals that contextually significant words disproportionately define phrases. This relationship is applied to produce simple weighted bag-of-words sentence embeddings. This model outperforms SkipThought and the best models trained on unordered sentences in most tests in Facebook's SentEval, beats tf-idf on all available tests, and is generally comparable to the state of the art. This paper also applies CoSal to sentence and document summarization and an improved and context-aware cosine distance. Applying the premise that unexpected words are important, CoSal is presented as a replacement for tf-idf and an intuitive measure of contextual word importance.

# 1 Introduction

# **Global Context**

Global context, a representation of the nature of text being analyzed, such as the recognition that one is reading a movie review or a sentence in a long news article about a specific topic, has long been acknowledged as a powerful tool for transfer learning and interpreting large documents [1][2]. However, current solutions to accounting for global context are essentially black boxes: the algorithm takes in a document and returns some function that takes a given word vector and returns its context measurement, which can be anything from a vector [3] to another deep structure [2]. This results in obfuscated approaches to context that are challenging to meaningfully analyze or invert and are resistant to further extension. Further, these algorithms often require thousands of sentences in a transfer dataset to train effectively, especially when unsupervised [3]. However, in this paper, using CoSal, a simple approach relying on normalization with respect to the word vector distribution of a context is developed, then used to augment several machine learning techniques.

It is generally accepted that the words that are common in a document but rare in a corpus are important: tf-idf is a standard technique to evaluate word importance, comparing the frequency of a word's use in a document to the frequency in the overall corpus, used by at least 70% of text based recommender-systems [4]. This premise also features in the inverse-frequency sentence embedding model in Arora et al.'s "A Simple but Tough to Beat Baseline for Sentence Embeddings" [5], which remains nearly state-of-the-art.

Thus, a natural extension is to measure the "rareness" of a word vector relative to its document's word vector distribution. Finding that CoSal correlates well with tf-idf, the natural follow-up analyzes how this correlates to the meaning of sentence, relative to the underlying words' importances. Building a vector space with both phrases and words led to a surprising but clear relationship: in essence, words that are slightly more contextually important than the others in a sentence contribute much more to the meaning of the sentence than words that are slightly less contextually important.

# **Sentence and Document Summarization**

Combining the CoSal pattern with a weighted bag-of-words yields unprecedented transfer learning performance. Since this linear combination of words should be in the same vector space as the initial words, one can the recursively extract the meaning from a sentence vector. For the algorithm, we assume one of a sentence vector's closest words is representative of the vector. Under this assumption, we develop a pathfinding algorithm from an empty word list to the word list best approximating the sentence embedding.

Initially, cosine distance was used to measure the distance between a generated sentence vector and a goal sentence vector, providing numerous advantages over Euclidean distance [6]. However, cosine distance implies that all dimensions are equally valuable and uncorrelated, as it is calculated as a dot product along the dimensions. Thus, it still produces noisy measurements. Noting that M-distance can be calculated between any two points with an underlying distribution, an application of the law of cosines gives rise to a more robust and context-aware alternative to cosine distance, capable of realizing that "cardinal" and "red" are less related to one another in the context of an essay about Stanford than in an article about the color green.

# 2 Background and Related Work

# 2.1 Sentence Embeddings in Context

The context of a dataset is often recognized as useful in tasks from interpreting a user's next search query to document summarization [1]. Further, the significance of global context was highlighted in "Deep contextualized word representations" [2] using a supervised model on a dataset to create a biLSTM encoder. In all of these papers, it is pointed out that the short-term nature of the memory of long short-term memory networks (LSTMs) stops such a model from learning global context, especially with short texts, and learning global contextual significance of a word requires degrees of freedom prone to over-fitting. Other research generated context-vectors for words, calculated by an autoencoder trained on a translation model [3] or used an approach to transform word vectors based on global context using a biLSTM [7]. While techniques have been developed for small scale transfer learning [1], the learning of global context [2], and multiple-word meaning word vectors [8], CoSal provides a more transparent, easy-to-implement technique which also works on tiny contexts.

# 2.2 tf-idf

tf-idf compares a word's frequency in a document to its frequency in an overall corpus to determine importance. tf-idf comes with a variety of issues. First, to produce meaningful results even for common words requires a large document. tf-idf is stratified for rarer tokens, resulting in importance bands for a document as words get increasingly rare. Further, it fails to provide useful results for out-of-document tokens, limiting its usefulness in comparisons. It also ignores covariance or similarity in word meaning, and thus is extremely word-choice dependent. A central idea of this paper is that tf-idf works not because it is good at detecting unusual words, but because the unusualness of words happens to be a good proxy for the unusualness of their meanings relative to the document/context.

# 2.3 Mahalanobis Distance

The Mahalanobis distance, or M-distance, is a simple way to find a distance between a point and a distribution or between two points sampled from a distribution, normalized for deviation and covariance [9]. The formula for two points  $p_1$  and  $p_2$  and distribution covariance S is:

$$d(p_1, p_2, S) = \sqrt{(p_1 - p_2)^T S^{-1} (p_1 - p_2)}$$

While the M-distance has been applied in the past in the context of word vectors, the uses have ranged from measuring the distance between "Gaussian word embeddings" <sup>1</sup> [10], to a substitute for cosine distance for one-shot image recognition <sup>2</sup> [11] to movie review sentiment analysis [12] or to help discriminate word sense

<sup>&</sup>lt;sup>1</sup>Word embeddings which have a mean and deviation in every dimension instead of simply a single point, for spanning ambiguities. The covariance used was the joint covariance of the word embeddings of the words in the context dataset

<sup>&</sup>lt;sup>2</sup>This technique results in rare words being treated as distant. For example "dog" and "and" are both common words, and will generally have a smaller M-distance than "Dachshund" and "dog."

in context [13]. Consistently, the measurement is used to measure the distance between two words in a sample, not the distance of a word to the distribution.

# 2.4 Statistical Approaches

As discussed in Arora et al.'s "A Simple but Tough to Beat Baseline for Sentence Embeddings" [5], a weighted average of words by their distance from the first principal component of a sentence yields a remarkably robust approximate sentence vector embedding. However, this "smooth inverse frequency" (SIF) approach comes with limitations. Not only is calculating PCA for every sentence in a document computationally complex, but the first principal component of a small number of normally distributed words in a high dimensional space is subject to random fluctuation. A variation to our sentence model draws from [5], incorporating the distance to the sentence average instead of the context average, using the context's covariance. Their calculation of word frequencies from the unigram count of the word in the corpus also means that their approach still does not work for out-of-vocab words, has no equivalent in other vector spaces and can't be generated from the word vectors alone, suffering some of the issues of tf-idf.

# 3 Analysis of M-Distance

# 3.1 Analyzing M-Distance vs tf-idf

One necessary initial test was an implementation of the Mahalanobis distance of words in a corpus compared to their tf-idf. The context used here was the Stanford Sentiment Analysis Treebank dataset [14]. Both normalized and unnormalized GloVe word vectors<sup>3</sup> [15] were compared, to see how much of the distance relationship was encoded in the initial distance of the word vectors. The normalized approach resulted in substantially better results, presumably because the length of the original word vectors was correlated with their frequency of use, which caused the document cloud to be less predictive.

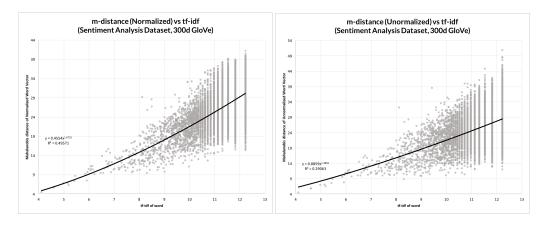


Figure 1: **Relationship of Mahalanobis distance to tf-idf** on Stanford Treebank Dataset with 300d GloVe word vectors, normalized (left) and unnormalized (right)

# 3.2 The Unified 2-gram and Token Space

Modifying an implementation of GloVe in TensorFlow, tf-glove [16], by using spaCy [17] to identify two-word phrase<sup>4</sup>, the algorithm randomly treated a two-word phrase as a single token 50% of the time, and the remaining times replaced it with one of the underlying words (To prevent the phrase's typical position in the sentence from biasing its word vector, given the reduced context at the beginning or end of a sentence). In order to analyze the relationship between the phrase vector and its constituent word vectors, the closest linear combination of the two words was used to produce the curve <sup>5</sup>.

<sup>&</sup>lt;sup>3</sup>300 dimensional word vectors from 42B token Common Crawl

<sup>&</sup>lt;sup>4</sup>Used in the grammatical sense of the word phrase

<sup>&</sup>lt;sup>5</sup>Note that the shape of the curve is almost identical when the phrases containing stop words are removed. Stop words can be approximately calculated as the least important 15% of words by distance

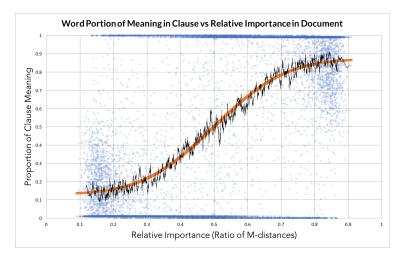


Figure 2: The proportion of a phrase's vector that is determined by a given word, as a function of the word's proportion of the sum of the two words' importances, shown with a moving average (black) and a sigmoid (orange, erf) fit to the data. The y-axis is calculated by finding the x such that  $x(v_1) + (1-x)(v_2)$  is as close to the phrase vector as possible. The x-axis is  $\frac{d(v_1)}{d(v_1)+d(v_2)}$ 

# 4 Broader Linguistic Context and CoSal

While the relationships in the previous section may make simply applying M-distance to the covariance as a measure of word importance appealing, that approach is limited: We hope to place less emphasis on features like tense or plurality that may be consistent in a text but vary in language or a larger contextual corpus.

For this, mirroring tf-idf, we weigh the covariance of the document and the covariance of the language or corpus, roughly corresponding to tf and idf respectively. Then, we take the (element-wise) product of the covariances. Letting S' be the weighted covariance S, CoSal is calculated as follows:

$$CoSal(v) = \sqrt{v^T(S'_{doc} \odot S'_{lang})^{-1}v}$$

Like tf-idf, CoSal supports arbitrary weighting schemes. The recommended weighting scheme, drawing from the common "augmented" [18] tf-idf weighting scheme<sup>6</sup>, with  $\sqrt[+]{x}$  as the sign-preserving square root

$$S_{corp}' = \sqrt[+]{S_{corp}}$$
 and  $S_{doc}' = \sqrt{rac{|S_{doc}| + |S_{corp}|}{2}}$ 

### Confidence

More generally,  $S'_{doc} = \sqrt{p|S_{doc}| + (1-p)|S_{corp}|}$ , where p is the proportion of information of the context, roughly ln(n)/10 where n is the context's word count. While multiplying the covariances without roots generally ranks words in the same order, this method maintains the shape of the distribution.

# 5 Sentence Embeddings

From this point in this paper onwards, GloVe is replaced with fastText<sup>7</sup> [19], which, has less predictive word vectors [20], but can perform a character level prediction to predict word vectors for misspelled or extremely esoteric words. This was used with PyMagnitude, which offered extremely fast word vector lookup [21].

# Challenges

A challenge of extending the two-word model to a sentence embedding model is that taking the sigmoid of the proportion of total importance will result in longer sentences having every word on the left side of the sigmoid. Thus, two times an average should replace the sum of the words. A variety of weights were tested, including a simple non-sigmoidal weight as a baseline, the peak-end approach which was computationally cheapest and surprisingly effective, and scaled by the mean, root-mean-square, and harmonic mean.

<sup>&</sup>lt;sup>6</sup>The "natural" [18] scheme, which is essentially just tf, is used in this paper for very large corpora

<sup>&</sup>lt;sup>7</sup>2 million word 300d vectors trained with Common Crawl, crawl-300d-2M

A recursive approach was also tested, where the sentence was broken down by phrase by spaCy [17] and then constituent phrases were recursively merged according to the sigmoid. Unfortunately, it clearly underperformed the bag-of-words approach from the start, so was dismissed.

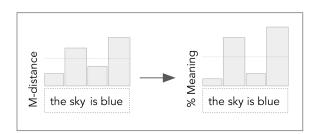


Figure 3: A visualization of the sigmoidal relationship between importance in a document and in a sentence.

# **Evaluation Technique**

Facebook's SentEval [20] suite of tests was used to evaluate the quality of the transfer-learned sentence embeddings, which includes subjectivity detection (SUBJ), product reviews (CR), entailment<sup>8</sup> (SICKEntailment), opinion polarity (MPQA), question-type identification (TREC), movie reviews (MR), paraphrase detection (MRPC), semantic similarity (STS), and other analyses. The tests perform a gradient-descent based classification on the returned sentence vectors.

Because the word vector model contained fewer capitalized words, case-insensitive performance was generally better than case-sensitive performance. With very small datasets (Many fewer words than dimensions), using lemmatization can improve performance.

#### Procedure

The following includes the variation discussed earlier in this paper as well as the global context only version, measuring to the sentence average instead of the global average. This variation approximates the sigmoid linearly for the local case, as the distances in the local case tend to cluster more. The global-only harmonic mean results are comparable.

**Data:** trainingSet, a large document or transfer learning training set

**Result:** Returns Mahalanobis metric of training set and average of training word vectors

vectorList =word vector of each word in trainingSet;

*qlobalCovariance*, *qlobalAverage* = covariance and average of *vectorList*;

**return** global Average and Mahalanobis metric using global Covariance,

Algorithm 1: Training algorithm. Finds the distribution of the document's word vector cloud

**Data:** sentence: a list of words; metric and globalAverage from trainingSet; globalOnly: a boolean for whether to use the global average or the sentence's average word vector.

**Result:** Weighs the words by their contextual importances

vecs = word vectors of sentence's words;

averageVec = globalAverage if globalOnly else normalize(average of vecs);

distanceList = metric distance between averageVec and each of vecs;

distanceList /=2\* average(distanceList);

weights = sigmoid(distanceList) if globalOnly else c \* (distanceList - 0.5) + 0.5;

return elementwise product of weights and vecs;

Algorithm 2: Sentence embedding algorithm. Returns the sentence vector in a trained context

# 5.1 Performance and Stability

This model requires orders of magnitude fewer examples than models like SkipThought or the bidirectional LSTM to produce high quality results and no training beyond that of the underlying word vectors. On a Macbook Pro, it takes 24ms per thousand sentences to "train" the model given training dataset word vectors. One

<sup>&</sup>lt;sup>8</sup>Whether a given sentence implies another sentence

question is how many samples are needed to establish a context. This model's performance exceeds tf-idf's for most tests after only 10 sentences. Further, every example seems to improve performance logarithmically. On the MRPC dataset, 99% of the model's ultimate accuracy is reached with 60 sentences, with 5, 10 and 30 example sentences resulting in accuracies of 94%, 96% and 97% of the ultimate performance.

Evaluation	tf-idf	ST	ST-LN	BiLSTM	GOBOW	SBOW
MR: Movie Review	73.7	76.5	79.4	77.5	78.4	78.0
MPQA: Opinion Polarity	82.4	87.1	89.3	88.7	87.7	88.3
MRPC: Paraphrase Detection	73.6/81.7	73.0/82.0	-	73.2/81.6	73.5/81.7	74.0/82.1
SICK-E: SICK Entailment	-	82.3	79.5	83.4	78.6	78.5
SST: Movie Review	-	82.0	82.9	80.7	82.6	82.8
SUBJ: Subjectivity/Objectivity	90.3	93.6	93.7	89.6	92.8	92.6
CR: Product Reviews Sentiment	79.2	80.1	83.1	81.3	81.5	80.8
TREC: Question Type	85.0	92.2	88.4	85.8	82.2	86.2

Table 1: Transfer Performance. The unsupervised state-of-the-art sentence embedding transfer learning technique, SkipThought with normalized layers, is listed as ST-LN, SkipThought, more widely used, is ST, A bidirectional LSTM was the best unordered unsupervised sentence model in [20]. Thus, BiLSTM-Max is used as a benchmark. Their results draw from [20]. The GOBOW model uses global context only, with a sigmoid through (0.5, 0.5) scaled vertically by 0.78, horizontally by 0.11 (slope = 1.80), from Figure 3, with the resulting vector normalized. The SBOW model, using the sentence average, was optimal at c = 1.87.

#### **Context-Adjusted Cosine Distance** 5.2

Due to the limitations of cosine distance, an approach that accounts for the underlying vector distribution was used repeatedly throughout the paper. Treating the CoSal of the difference between the words as the opposite leg of the triangle and treating the word's salience measures as the legs, the law of cosines is solved for the cosine of the angle between the words 9. For example, in the context of an article about a development in Stanford self-driving car tech, "cardinal" and "red" have a distance of 0.94312, while in random excerpts from a Wikipedia article about the color green<sup>10</sup> [22], they have a distance of 0.908867.

#### 5.3 Sentence Summarization

Since these sentence embeddings are in the same space as the word embeddings, one of a sentence vector's

closest words reflects its meaning. By the sigmoidal phrase sum relationship, where d is the M-distance, 
$$v_{sent} = \sigma(\frac{d(v_1)}{d(v_1) + d(v_2)})v_1 + \sigma(\frac{d(v_2)}{d(v_1) + d(v_2)})v_2, \text{ it is also true that } \frac{v_{sent} - \sigma(\frac{d(v_1)}{d(v_1) + d(v_2)})v_1}{\sigma(\frac{d(v_2)}{d(v_1) + d(v_2)})} = v_2. \text{ Starting with } v_1 = v_2$$

 $d(v_2) = d_{avg}$  and calculating  $v_2$ , the resulting  $d(v_2)$  can be plugged back in, repeatedly. This converges quickly and allows the calculation of the vector of a sentence vector without some words. Thus, in order to summarize a sentence from the sentence vector, the closest 5 words to the remaining sentence are chosen. Then, they are fed into a priority queue by context-adjusted cosine distance until the target sentence vector is within a given radius. This slightly modified A\* algorithm extracts meaning from a sentence vector.

- (0.0) spiderman rocks  $\Rightarrow$  spiderman rocks
- (0.304813046595) light, cute and forgettable  $\Rightarrow$  trifle believable glowing forgettable cute light cuter
- (0.208039930017) not so much farcical as sour.  $\Rightarrow$  laughable so much nothing sour farcical
- (0.35074034786) effective but too-tepid biopic  $\Rightarrow$  tepid know-but slug-fest precise biopic effective
- (0.270376983591) occasionally melodramatic, it 's also extremely effective. ⇒ unbelievably sometimes melodramatically effective though occasionally
- (0.366886138345) a visually flashy but narratively opaque and emotionally vapid exercise in style and mystification. ⇒ theatricality verisimilitude dull mystifying bland opaque narratively

Figure 4: A generated summarization for some of the movie review test sentences. Calculated with the fastest heuristic, requiring every next word to reduce the distance by at least 5%.

 $<sup>\</sup>frac{9}{2}\frac{a^2+b^2-c^2}{2ab}$  where a and b are the legs and c is the opposite side  $^{10}$ As appearing on March 18

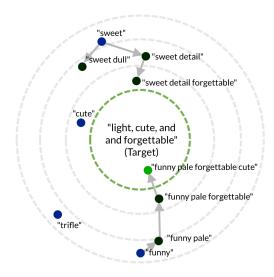


Figure 5: A pathfinding visualization. Each next word is one of the words closest to the remaining vector.

# 5.4 Example Application of Sigmoid CoSal

Decid	0.60	Hz.	0.00	III.	0.00	rades.	0.00	N A le	0.00	0	0.54	11.	0.00	restant	0.70
	0.62		0.29		0.28		0.68			Overall	0.54			Erdoss	0.79
	0.73					J	0.71	of		his		firmly		prolific	0.69
			0.40		0.34		0.20			work		believed		output	0.70
						January Inches	0.71			leaned		mathematics		with	0.32
						problems	0.60	centered		towards	0.58			coauthors	0.77
			0.35		0.28	in	0.19	around		solving	0.54			prompted	0.80
				mathematics			0.60	discrete		previously	0.57			the	0.30
1996			0.63	even	0.24	mathematics	0.56		0.64		0.61	social		creation	0.55
was	0.30				0.41	graph	0.73	cracking		problems	0.62	activity	0.68	of	0.20
a	0.22	mathematics	0.68	his	0.34	theory	0.45	many	0.27	rather	0.29	living	0.81		0.30
Hungarian	0.75	he			0.30		0.38	previously	0.63	than	0.30		0.45	Erdos	0.76
mathematician	0.68	engaged	0.81	yearsindeed	0.79	theory	0.45	unsolved	0.82	developing	0.68	itinerant	0.75	number	0.46
He	0.37	more	0.38	his	0.34	mathematical	0.52	problems	0.68	or	0.34	lifestyle	0.69	the	0.30
was	0.43	than	0.35	death	0.56	analysis	0.71	in	0.21	exploring	0.74	with	0.30	number	0.46
one	0.31	500	0.70	came	0.35	approximation	0.68	the	0.29	new	0.35	the	0.28	of	0.20
of	0.23	collaborators	0.79	only	0.73	theory	0.45	field	0.60	areas	0.77	sole	0.71	steps	0.71
the	0.39	and	0.23	hours	0.76	set	0.35	He	0.27	of	0.18	purpose	0.60	in	0.21
most	0.41	for	0.26	after	0.35	theory	0.45	championed	0.79	mathematics	0.58	of	0.19	the	0.30
prolific	0.80	his	0.35	he	0.28	and	0.20	and	0.22	Erdos	0.68	writing	0.55	shortest	0.78
mathematicians	0.79	eccentric	0.74	solved	0.83	probability	0.60	contributed	0.77	published	0.58	mathematical	0.58	path	0.70
of	0.23	lifestyle	0.74	a	0.24	theory	0.45	to	0.26	around	0.41	papers	0.73	between	0.52
the	0.39	Time	0.29	geometry	0.75			Ramsey	0.81	1500	0.74	with	0.30	a	0.23
20th	0.76	magazine	0.67	problem	0.56			theory	0.52	mathematical	0.51	other	0.28	mathematician	0.71
century	0.76	called	0.77	in	0.21			which	0.34	papers	0.67	mathematicians	0.64	and	0.22
,		him	0.39	a	0.24			studied	0.75	during	0.34			Erdos	0.76
		The	0.42	conference	0.57			the	0.29	his	0.27			in	0.21
		Oddballs	0.74	in	0.21			conditions	0.72	lifetime	0.66			terms	0.75
		Oddball	0.78	Warsaw	0.78			in	0.21	a	0.21			of	0.20
	,					*		which	0.34	figure	0.59			coauthorships	0.74
								order	0.52	· ·	0.21				
								necessarily	0.67	remains	0.51				
				appears	0.47	unsurpassed	0.73								

Figure 6: The global sigmoid weights generated for a short excerpt about Erdos from Wikipedia [23], using only the shown text as document context and TREC for linguistic context, using the recommended weighting with p=0.2

# 5.5 Document Summarization

Applying the same basic technique used for sentence summarization to documents, the sentence vector of every sentence is calculated, then the document vector is calculated as a composition of sentence vectors using the CoSal of the sentences. Notably, treating the document as a sentence produces qualitatively worse summarizations. The context-adjusted cosine distance is then used to find the closest sentences, then to remove those sentence vectors from the document meaning, until the best distance for a given depth is reached.

(0.022 distance) A driverless car is making its way through a winding neighborhood street, about to make a sharp turn onto a road where a child's ball has just rolled. Although no person in the car can see that ball, the car stops to avoid it ... 622 words ... In the future, we want to make it even more practical in the 'wild.

⇒ Before this system is road ready, it will also have to work better in daylight and with objects in motion, like a bouncing ball or running child.

================

(0.091 distance) In surveys made in American, European, and Islamic countries, green is the color most commonly associated with nature, life, health, youth, spring, hope and envy. In the European Union and the United States, green is also sometimes associated with toxicity and poor health, ... 579 words ... to create a dark blue sky with a greenish tint in his painting Cafe terrace at night

⇒ In more contemporary terms, they are. guriin, which is derived from the English word green). green. Because of its association with nature, it is the color of the environmental movement.

(0.10 distance) Depending on your background you might be wondering: What makes Recurrent Networks so special? A glaring limitation of Vanilla Neural Networks (and also Convolutional Networks) is that their API is too constrained: they accept a fixed-sized vector as input (e.g. an image) and produce a fixed-sized vector ... 474 words ... On the right, a recurrent network generates images of digits by learning to sequentially add color to a canvas (Gregor et al.)

⇒ Vanilla mode of processing without RNN, from fixed-sized input to fixed-sized output. Sequences in the input, the output, or in the most general case both. This can in programming terms be interpreted as running a fixed program with certain inputs and some internal variables. Notice that in every case are no pre-specified constraints on the lengths sequences because the recurrent transformation.

Figure 7: Generated summarizations of Stanford News' "Stanford researchers develop technique to see objects hidden around corners" [24]; random excerpts from the Wikipedia article on green (From March 18, from beginning to "Languages where green and blue are one color," and "Pigments," and "Color vision") [22]; an extract of Karpathy's famous "Unreasonable Effectiveness of Recurrent Neural Networks" [25]

# 6 Conclusion and Future Directions

Beyond useful calculations to augment linguistic algorithms, a novel summarization technique, and a valuable new metric, this paper carries implications about the compositional structure of language. Fundamentally, the idea that context allows people to consistently disregard and overemphasize information when analyzing collective meanings is profound. It suggests that, although all parts play a role in establishing collective meaning, for example, of a series of novels, what comes to mind when the series is thought of is disproportionately a few books that stand out. The recursive document embedding pattern suggests consolidation of information can aid abstraction, possibly due to working memory limits. This may provide insight into the neurological foundations of higher-order reasoning.

CoSal is a natural and broadly applicable technique for contextual importance. It can be performed on any data representable in a vector space with mostly normally distributed dimensions and distance corresponding to similarity. For example, CoSal should be able to identify important frames in a video or songs in a playlist as image and song vectors. Additionally, the base linguistic salience technique can likely be augmented with contextual word vectors [2]. Ideally, especially given the ease of implementation, CoSal should replace tf-idf.

# Acknowledgments

I would like to thank Nick Cammarata for his substantial and useful feedback on this paper, its presentation, and for noting that Contextual Salience is generalizable beyond text. I would also like to thank Sarina Wu, Samuel Jacobo, and Christian Cosgrove for reading over this paper and making suggestions for clarity. Finally, I would like to express thanks to Zewei Chu, who suggested using SentEval, with whom I ran into the limitations of tf-idf and its associated sentence embedding approach.

# References

- [1] Z. Lu, Y. Zhu, S. Pan, E. Xiang, Y. Wang, and Q. Yang, "Source free transfer learning for text classification," 2014.
- [2] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," *ArXiv e-prints*, Feb. 2018.
- [3] B. McCann, J. Bradbury, C. Xiong, and R. Socher, "Learned in Translation: Contextualized Word Vectors," *ArXiv e-prints*, July 2017.
- [4] J. Beel, B. Gipp, S. Langer, and C. Breitinger, "Research-paper recommender systems: a literature survey," *International Journal on Digital Libraries*, vol. 17, pp. 305–338, Nov 2016.
- [5] S. Arora, Y. Liang, and T. Ma, "A simple but tough-to-beat baseline for sentence embeddings," 2017.
- [6] C. Emmery, "Euclidean vs. cosine distance," Mar 2017.
- [7] M. E. Peters, W. Ammar, C. Bhagavatula, and R. Power, "Semi-supervised sequence tagging with bidirectional language models," *ArXiv e-prints*, Apr. 2017.
- [8] E. H. Huang, R. Socher, C. D. Manning, and A. Y. Ng, "Improving word representations via global context and multiple word prototypes," in *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2012.
- [9] R. Gnanadesikan and J. R. Kettenring, "Robust estimates, residuals, and outlier detection with multiresponse data," *Biometrics*, vol. 28, no. 1, pp. 81–124, 1972.
- [10] L. Vilnis and A. McCallum, "Word Representations via Gaussian Embedding," ArXiv e-prints, Dec. 2014.
- [11] S. Rahman, S. H. Khan, and F. Porikli, "A Unified approach for Conventional Zero-shot, Generalized Zero-shot and Few-shot Learning," *ArXiv e-prints*, June 2017.
- [12] V. Balasubramanian, S. Gupta, and P. Veerappagoundar, "Mahalanobis distance-the ultimate measure for sentiment analysis," vol. 13, pp. 252–257, 01 2016.
- [13] M.-C. De Marneffe, C. Archambeau, P. Dupont, and M. Verleysen, "Local vector-based models for sense discrimination," 03 2018.
- [14] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," 2013.
- [15] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation.," in *EMNLP*, vol. 14, pp. 1532–1543, 2014.
- [16] G. Simon, "tensorflow-glove," Mar 2017.
- [17] M. Honnibal and M. Johnson, "An improved non-monotonic transition system for dependency parsing," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, (Lisbon, Portugal), pp. 1373–1378, Association for Computational Linguistics, September 2015.
- [18] C. D. Manning, P. Raghavan, and S. Hinrich, *Introduction to information retrieval*. Cambridge University Press, 2009.
- [19] T. Mikolov, E. Grave, P. Bojanowski, C. Puhrsch, and A. Joulin, "Advances in pre-training distributed word representations," in *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [20] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, and A. Bordes, "Supervised learning of universal sentence representations from natural language inference data," *arXiv preprint arXiv:1705.02364*, 2017.
- [21] A. Patel and Alex, "plasticityai/magnitude: Release 0.1.14," Mar. 2018.
- [22] "Green," Mar 2018.
- [23] "Paul Erds," Mar 2018.
- [24] T. Kubota, "Technique can see objects hidden around corners," Mar 2018.
- [25] A. Karpathy, "The unreasonable effectiveness of recurrent neural networks," May 2015.