

Appendix A Model Hyperparameters

Similar to Jain and Wallace (2019) we use Fast-Text pretrained embeddings (Joulin et al., 2016) for the SST and ADR datasets, Glove pretrained embeddings (Pennington et al., 2014) for the IMDB and AG News datasets, while we use Word2Vec (Mikolov et al., 2013) from Gensim (Řehůřek and Sojka, 2010) to train embeddings for MIMIC. All embeddings are of size $d = 300$. We also replace all numbers in text with a special symbol q and initialise the embeddings of unknown words randomly from a normal distribution, $\mathcal{N}(0, 1)$. The embeddings are not trained alongside the rest of the model.

We train the models using default Adam learning rate ($1e-3$) with $1e-4$ weight decay, which adds an l_2 regulariser across all parameters. We use 64 dimensional hidden representations for one-layered bi-LSTM and bi-GRU encoders and 128 dimensional hidden representation for the MLP encoder following Jain and Wallace (2019). For the CNN we use 4 kernels of sizes $[1, 3, 5, 7]$, each with 32 filters, giving a final contextual representation \mathbf{h}_i of size $N = 128$, with ReLU activation function on the output of the filters, as per Jain and Wallace (2019).

For BERT we use the pre-trained version from Wolf et al. (2020) and fine-tune with a learning rate of $1e - 5$ all BERT parameters except from the word embeddings, to simulate the scenario with the rest of the encoders, and $1e - 4$ for the remainder of the parameters. We train our models three times using different random seeds and a batch size of 8 for BERT and 32 for the rest of the models.

For Conv-TaSc we apply a CNN with 15 channels over the scaled embedding e_i from Lin-TaSc, keeping a single stride and a 1-dimensional kernel. This way, we ensure that input words remain context-independent. We then sum over the filtered scaled embedding e_i^f , to obtain the scores s_{x_i} . We have also experimented with filter sizes of $[2, 10, 20, 30, 50]$ individually and simultaneously.

For the MIMIC dataset we also attempted to use LongFormer (Beltagy et al., 2020), which is a BERT version that has the ability to accept and deal with longer sequences. However due to the increasing time to train and evaluate the model, this BERT variant was abandoned. Additionally we attempted to use Hierarchical BERT to deal with the longer sequences, however increases where not substantial and run times where similarly increased. Finally,

contrary to the remainder of the datasets to deal with the long sequences of MIMIC we truncated the 256 first tokens and 256 last tokens, following the suggestions of Sun et al. (2019). We experimented with using the first and the last 512 tokens, but the head and tails truncation approach yielded the best performances.

Appendix B Additional parameters with TaSc variants

In Table 6 we present the additional parameters introduced by each variant, with Lin-TaSc requiring the lowest number of parameters and Feat-TaSc the most.

TaSc Mechanism	Additional Parameters
Lin-TaSc	$ \mathbf{V} $
Feat-TaSc	$ \mathbf{V} \times d$
Conv-TaSc	$ \mathbf{V} + d \times n + n$

Table 6: Additional parameters resulting from the proposed TaSc mechanisms where $|\mathbf{V}|$ is the vocabulary size, d the embedding dimension and n the number of channels in a CNN.

Appendix C Reproducibility Results

Computational infrastructure used: For the experiments above we used NVIDIA’s TESLA V100 GPU.

Dataset description: We consider the following datasets for text classification following Wiegrefe and Pinter (2019) and Jain and Wallace (2019):

SST: *Stanford Sentiment Treebank* consists of sentences tagged with sentiment on a 5-point-scale from negative to positive (Socher et al., 2013). Jain and Wallace (2019) removed sentences with neutral sentiment and labelled the remaining sentences to negative and positive if they have a score lower or higher than 3 respectively.

IMDB: The *Large Movie Reviews Corpus* consists of 50,000 movie reviews labelled either as positive or negative (Maas et al., 2011). We filter the dataset as per Jain and Wallace (2019) to include movie reviews with sequence length less than 400 words.

ADR: A dataset of $\sim 20,000$ tweets with labels indicating whether a Twitter post contains an adverse drug reaction or not (Sarker et al., 2015).

Data-set	Enc()	No-TaSc		Lin-TaSc		Feat-TaSc		Conv-TaSc	
		Dot	Tanh	Dot	Tanh	Dot	Tanh	Dot	Tanh
SST	BERT	.89	.90	.90	.87	.87	.87	.90	.90
	LSTM	.77	.78	.77	.78	.77	.80	.79	.80
	GRU	.78	.78	.78	.79	.78	.79	.78	.79
	MLP	.75	.77	.78	.78	.80	.80	.79	.81
	CNN	.77	.77	.79	.80	.80	.79	.79	.78
ADR	BERT	.81	.81	.81	.79	.80	.80	.80	.81
	LSTM	.74	.75	.77	.76	.77	.77	.78	.76
	GRU	.76	.75	.77	.77	.76	.79	.77	.77
	MLP	.73	.78	.76	.76	.78	.77	.76	.76
	CNN	.74	.73	.77	.76	.77	.77	.78	.78
IMDB	BERT	.92	.92	.93	.92	.92	.92	.92	.92
	LSTM	.90	.89	.89	.89	.89	.89	.89	.89
	GRU	.90	.90	.89	.90	.89	.90	.89	.89
	MLP	.88	.88	.88	.88	.89	.88	.89	.88
	CNN	.89	.89	.90	.89	.89	.89	.89	.89
AG	BERT	.95	.95	.94	.94	.95	.95	.94	.95
	LSTM	.93	.93	.92	.93	.93	.93	.93	.93
	GRU	.93	.93	.93	.93	.93	.93	.93	.93
	MLP	.93	.93	.93	.92	.93	.92	.93	.93
	CNN	.93	.93	.93	.93	.93	.93	.93	.93
MIMIC	BERT	.84	.83	.85	.84	.86	.84	.85	.83
	LSTM	.88	.89	.89	.89	.89	.90	.90	.90
	GRU	.89	.90	.89	.89	.90	.90	.90	.90
	MLP	.90	.89	.88	.88	.89	.88	.89	.89
	CNN	.90	.89	.90	.90	.90	.90	.89	.90

Table 7: Validation set F1-macro average scores (3 runs) across datasets, encoders and attention mechanisms for models with and without TaSc (No-TaSc). Standard deviations do not exceed 0.01.

AG: A subset of the original news articles¹³ dataset compiled by [Jain and Wallace \(2019\)](#) for topic categorisation (*Business* and *World* news).

MIMIC: A sample of discharge summaries from the MIMIC III dataset of health records ([Johnson et al., 2016](#)). The task is to recognise if a given summary has been labelled as relevant to acute or chronic anemia ([Jain and Wallace, 2019](#)).

Validation set predictive performances: In Table 7 we present predictive performances on the validation checks for reproducibility on models with TaSc and models without (No-TaSc).

Appendix D Detailed Experiment Results

In Tables 8, 9, 10 and 11 we present results in numbers for Figures 1 and 2. Tables 8 and 9 show the percentage of decision flips recorded by removing the highest scored token, across encoders and datasets respectively and compliment Figure 1. Tables 10 and 11 show the average fraction of tokens required to be masked in order to cause a change in

¹³https://di.unipi.it/~gulli/AG_corpus_of_news_articles.html. Accessed on Sep 2019

	Enc()	No-TaSc		Lin-TaSc		Feat-TaSc		Conv-TaSc	
		Dot	Tanh	Dot	Tanh	Dot	Tanh	Dot	Tanh
α	BERT	4.8		6.2 (1.3)		6.6 (1.4)		3.7 (0.8)	
	LSTM	6.2		4.8 (0.8)		5.1 (0.8)		4.8 (0.8)	
	GRU	6.2		5.7 (0.9)		5.5 (0.9)		5.4 (0.9)	
	MLP	8.0		6.0 (0.8)		5.2 (0.7)		5.6 (0.7)	
	CNN	9.3		6.2 (0.7)		5.7 (0.6)		5.4 (0.6)	
$\nabla\alpha$	BERT	5.2		6.4 (1.2)		7.4 (1.4)		3.6 (0.7)	
	LSTM	6.5		10.9 (1.7)		12.5 (1.9)		12.0 (1.9)	
	GRU	6.3		11.3 (1.8)		12.4 (2.0)		11.8 (1.9)	
	MLP	9.8		12.0 (1.2)		13.1 (1.3)		13.0 (1.3)	
	CNN	10.1		12.0 (1.2)		13.2 (1.3)		13.4 (1.3)	
$\alpha\nabla\alpha$	BERT	5.7		8.0 (1.4)		9.3 (1.6)		4.0 (0.7)	
	LSTM	8.3		12.8 (1.5)		13.6 (1.6)		13.1 (1.6)	
	GRU	8.3		14.2 (1.7)		13.9 (1.7)		13.4 (1.6)	
	MLP	13.7		14.6 (1.1)		13.9 (1.0)		13.8 (1.0)	
	CNN	13.9		14.7 (1.1)		14.5 (1.0)		14.6 (1.0)	

Table 8: Mean average *percentage of decision flips* occurred by removing the most informative token, using the three TaSc variants and No-TaSc across encoders (higher is better).

	Dataset	No-TaSc		Lin-TaSc		Feat-TaSc		Conv-TaSc	
		Dot	Tanh	Dot	Tanh	Dot	Tanh	Dot	Tanh
α	SST	16.7		12.5 (0.7)		11.2 (0.7)		11.5 (0.7)	
	ADR	4.5		5.8 (1.3)		5.4 (1.2)		3.9 (0.9)	
	IMDB	6.6		5.0 (0.8)		5.2 (0.8)		3.9 (0.6)	
	AG	3.5		3.1 (0.9)		3.8 (1.1)		2.7 (0.8)	
	MIMIC	3.0		2.5 (0.8)		2.5 (0.8)		2.5 (0.8)	
$\nabla\alpha$	SST	18.8		25.8 (1.4)		27.5 (1.5)		25.4 (1.3)	
	ADR	5.2		9.4 (1.8)		10.6 (2.0)		8.1 (1.5)	
	IMDB	7.1		7.6 (1.1)		9.2 (1.3)		9.3 (1.3)	
	AG	4.3		5.1 (1.2)		6.4 (1.5)		6.4 (1.5)	
	MIMIC	2.5		4.8 (1.9)		4.8 (1.9)		4.8 (1.9)	
$\alpha\nabla\alpha$	SST	24.1		29.4 (1.2)		29.5 (1.2)		27.2 (1.1)	
	ADR	6.0		10.3 (1.7)		11.1 (1.8)		8.2 (1.4)	
	IMDB	10.3		12.0 (1.2)		11.0 (1.1)		10.8 (1.0)	
	AG	5.2		6.4 (1.2)		8.0 (1.5)		6.9 (1.3)	
	MIMIC	4.3		6.2 (1.5)		5.7 (1.3)		5.8 (1.3)	

Table 9: Mean average *percentage of decision flips* occurred by removing the most informative token, using the three TaSc variants and No-TaSc across datasets (higher is better).

	Enc()	No-TaSc		Lin-TaSc		Feat-TaSc		Conv-TaSc	
		Dot	Tanh	Dot	Tanh	Dot	Tanh	Dot	Tanh
α	BERT	.59		.46 (0.8)		.44 (0.7)		.56 (0.9)	
	LSTM	.56		.48 (0.9)		.51 (0.9)		.52 (0.9)	
	GRU	.57		.45 (0.8)		.49 (0.9)		.50 (0.9)	
	MLP	.41		.43 (1.0)		.44 (1.1)		.46 (1.1)	
	CNN	.45		.47 (1.0)		.47 (1.0)		.44 (1.0)	
$\nabla\alpha$	BERT	.52		.34 (0.6)		.31 (0.6)		.58 (1.1)	
	LSTM	.44		.21 (0.5)		.17 (0.4)		.19 (0.4)	
	GRU	.46		.17 (0.4)		.18 (0.4)		.19 (0.4)	
	MLP	.21		.18 (0.8)		.17 (0.8)		.17 (0.8)	
	CNN	.30		.17 (0.6)		.17 (0.6)		.17 (0.6)	
$\alpha\nabla\alpha$	BERT	.52		.29 (0.6)		.29 (0.6)		.57 (1.1)	
	LSTM	.44		.20 (0.5)		.16 (0.4)		.18 (0.4)	
	GRU	.43		.16 (0.4)		.17 (0.4)		.18 (0.4)	
	MLP	.17		.15 (0.9)		.16 (0.9)		.16 (0.9)	
	CNN	.27		.16 (0.6)		.16 (0.6)		.16 (0.6)	

Table 10: Mean average *fraction of tokens* required to cause a decision flip, using the three TaSc variants and No-TaSc across encoders (lower is better).

prediction (decision flip) and compliment Figure 2.

	Dataset	No-TaSc	Lin-TaSc	Feat-TaSc	Conv-TaSc
α	SST	.45	.49 (1.1)	.48 (1.1)	.50 (1.1)
	ADR	.88	.72 (0.8)	.76 (0.9)	.77 (0.9)
	IMDB	.38	.31 (0.8)	.36 (1.0)	.41 (1.1)
	AG	.59	.52 (0.9)	.49 (0.8)	.54 (0.9)
	MIMIC	.28	.24 (0.9)	.26 (0.9)	.26 (0.9)
$\nabla\alpha$	SST	.35	.24 (0.7)	.20 (0.6)	.27 (0.8)
	ADR	.78	.38 (0.5)	.36 (0.5)	.47 (0.6)
	IMDB	.19	.10 (0.5)	.12 (0.6)	.14 (0.7)
	AG	.49	.35 (0.7)	.28 (0.6)	.35 (0.7)
	MIMIC	.13	.02 (0.2)	.03 (0.3)	.07 (0.5)
$\alpha\nabla\alpha$	SST	.33	.23 (0.7)	.19 (0.6)	.26 (0.8)
	ADR	.77	.34 (0.5)	.35 (0.5)	.47 (0.6)
	IMDB	.17	.07 (0.4)	.10 (0.6)	.13 (0.7)
	AG	.46	.31 (0.7)	.25 (0.5)	.34 (0.7)
	MIMIC	.10	.02 (0.2)	.03 (0.3)	.06 (0.7)

Table 11: Mean average *fraction of tokens* required to cause a decision flip, using the three TaSc variants and No-TaSc across datasets (lower is better).

Appendix E Comparing TaSc with Non-attention Input Importance Metrics

Data	Enc()	Tanh							Dot						
		Non - TaSc			Lin-TaSc				Non-TaSc			Lin-TaSc			
		WO	$x\nabla x$	IG	WO	$x\nabla x$	IG	$\alpha\nabla\alpha$	WO	$x\nabla x$	IG	WO	$x\nabla x$	IG	$\alpha\nabla\alpha$
SST	BERT	.29	.64	.51	.37	<u>.30</u>	<u>.25</u>	.22	.32	.62	.49	.35	<u>.57</u>	.51	.55
	LSTM	.25	.24	.20	.26	.33	.19	.19	.21	.23	.19	.21	.19	.19	.19
	GRU	.24	.22	.19	.29	.24	.20	.18	.24	.25	.23	<u>.21</u>	.19	.19	.19
	MLP	.36	.26	.24	<u>.26</u>	<u>.20</u>	<u>.19</u>	.18	.22	.19	.18	.24	.19	.19	.18
	CNN	.30	.25	.20	<u>.27</u>	<u>.22</u>	<u>.20</u>	.19	.22	.20	.18	<u>.21</u>	<u>.20</u>	<u>.20</u>	.19
ADR	BERT	.83	.91	.89	<u>.73</u>	<u>.55</u>	<u>.38</u>	.31	.81	.90	.87	<u>.68</u>	<u>.58</u>	<u>.52</u>	.50
	LSTM	.82	.81	.80	<u>.54</u>	<u>.42</u>	<u>.34</u>	.32	.87	.88	.87	<u>.42</u>	<u>.35</u>	.34	.34
	GRU	.84	.84	.84	<u>.49</u>	<u>.38</u>	<u>.36</u>	.35	.79	.80	.80	<u>.50</u>	<u>.40</u>	<u>.44</u>	.38
	MLP	.71	.63	.57	<u>.60</u>	<u>.36</u>	<u>.43</u>	.31	.49	.43	.39	<u>.49</u>	<u>.40</u>	<u>.44</u>	.40
	CNN	.80	.78	.78	<u>.57</u>	<u>.46</u>	<u>.39</u>	.37	.77	.74	.74	<u>.52</u>	<u>.43</u>	<u>.38</u>	.36
IMDB	BERT	.23	.69	.43	.27	.14	.16	.07	.24	.72	.49	.26	<u>.27</u>	.17	<u>.20</u>
	LSTM	.18	.12	.07	<u>.11</u>	.13	<u>.05</u>	.04	.26	.09	.07	<u>.07</u>	<u>.06</u>	<u>.06</u>	.05
	GRU	.18	.12	.07	<u>.11</u>	<u>.06</u>	<u>.05</u>	.04	.27	.15	.08	<u>.09</u>	.05	.05	.05
	MLP	.16	.05	.05	<u>.07</u>	.05	.05	.05	.18	.07	.06	<u>.09</u>	.05	.05	.05
	CNN	.21	.09	.07	<u>.18</u>	<u>.07</u>	<u>.06</u>	.05	.27	.07	.06	<u>.14</u>	<u>.07</u>	<u>.06</u>	.05
AG	BERT	.62	.78	.56	.64	<u>.58</u>	<u>.54</u>	.50	.56	.76	.60	.56	<u>.59</u>	.55	<u>.60</u>
	LSTM	.53	.51	.30	<u>.47</u>	<u>.37</u>	.31	.38	.47	.52	.35	<u>.43</u>	<u>.40</u>	.36	.46
	GRU	.45	.36	.31	.50	<u>.30</u>	<u>.24</u>	.20	.54	.40	.30	<u>.36</u>	<u>.24</u>	<u>.23</u>	.22
	MLP	.53	.24	.25	.53	<u>.23</u>	<u>.23</u>	.19	.44	.25	.23	<u>.40</u>	.19	<u>.25</u>	.19
	CNN	.55	.38	.28	<u>.48</u>	<u>.29</u>	<u>.24</u>	.20	.53	.35	.25	<u>.39</u>	<u>.27</u>	<u>.23</u>	.21
MIMIC	BERT	.24	.67	.43	.31	<u>.10</u>	<u>.04</u>	.03	.21	.57	.26	.25	<u>.07</u>	.05	.05
	LSTM	.35	.32	.12	.37	<u>.01</u>	<u>.02</u>	.01	.28	.40	.30	.40	<u>.01</u>	<u>.02</u>	.01
	GRU	.20	.24	.23	.46	<u>.01</u>	<u>.02</u>	.01	.36	.18	.08	.42	<u>.01</u>	<u>.02</u>	.01
	MLP	.40	.03	.22	<u>.18</u>	<u>.01</u>	<u>.02</u>	.01	.13	.04	.03	.16	<u>.02</u>	<u>.02</u>	.02
	CNN	.26	.15	.02	.52	<u>.01</u>	<u>.01</u>	.01	.43	.09	.02	.49	<u>.03</u>	<u>.02</u>	.02

Table 12: Average fraction of tokens required to cause a decision flip using the best performing attention-based ranking ($\alpha\nabla\alpha$) with TaSc, *Word omission*, (**WO**), *InputXGrad*, (∇x) and *Integrated Gradients* (**IG**). Underlined values denote that Lin-TaSc is better and **bold** values denote the best performing method row-wise. (lower is better)