

cases	doc_1		doc_2		decision	id
					DUPLICATES	189
			authors	<ul style="list-style-type: none">Katarína BenešováAndrej Á vecMarek Á uppa		
	authors	<ul style="list-style-type: none">Katarína BenešováAndrej Á vecMarek Á uppa	title	Cost-effective Deployment of BERT Models in Serverless Environment		
	title	Cost-effective Deployment of BERT Models in Serverless Environment	publication_date	2021-04-19 00:00:00		
	publication_date	2021-03-19 00:00:00	source	SupportedSources.INTERNET_ARCHIVE		
	source	SupportedSources.OPENALEX	journal			
	journal	arXiv (Cornell University)	volume			
	volume		doi			
	doi	None	urls	<ul style="list-style-type: none">https://web.archive.org/web/20210421032539/https://arxiv.org/pdf/2103.10673v2.pdf		
	urls	<ul style="list-style-type: none">https://openalex.org/W3137254616	id	id-5406197705243831488		
	id	id-1879542727168961955	abstract	In this study we demonstrate the viability of deploying BERT-style models to serverless environments in a production setting. Since the freely available pre-trained models are too large to be deployed in this way, we utilize knowledge distillation and fine-tune the models on proprietary datasets for two real-world tasks: sentiment analysis and semantic textual similarity. As a result, we obtain models that are tuned for a specific domain and deployable in serverless environments. The subsequent performance analysis shows that this solution results in latency levels acceptable for production use and that it is also a cost-effective approach for small-to-medium size deployments of BERT models, all without any infrastructure overhead.		
	abstract		versions			
	versions					