

cases	doc_1		doc_2				decision	id
							DUPLICATES	187
	authors	<ul style="list-style-type: none"><li>KarolĀna BeneĀovĀĵ</li><li>Andrej Ā vec</li><li>Marek Ā uppa</li></ul>	authors	<ul style="list-style-type: none"><li>Katar'ina Benevsov'a</li><li>Andrej vSvec</li><li>Marek vSuppa</li></ul>				
	title	Cost-effective Deployment of BERT Models in Serverless Environment	title	Cost-effective Deployment of BERT Models in Serverless Environment				
	publication_date	2021-03-19 00:00:00	publication_date	2021-03-19 00:00:00				
	source	SupportedSources.OPENALEX	source	SupportedSources.SEMANTIC_SCHOLAR				
	journal	arXiv (Cornell University)	journal					
	volume		volume					
	doi	None	doi					
	urls	<ul style="list-style-type: none"><li>https://openalex.org/W3137254616</li></ul>	urls	<ul style="list-style-type: none"><li>https://www.semanticscholar.org/paper/2101193a3ec4d8fe260c7505614e760a7235ecf5</li></ul>				
	id	id-1879542727168961955	id	id8821586626506311802				
	abstract		abstract	In this study we demonstrate the viability of deploying BERT-style models to AWS Lambda in a production environment. 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 the serverless environment. The subsequent performance analysis shows that this solution does not only report latency levels acceptable for production use but that it is also a costeffective alternative to small-to-medium size deployments of BERT models, all without any infrastructure overhead.				
	versions		versions					