

501 **G Dataset Information**

502 **Dataset Card.** The dataset card is included on this page [https://huggingface.co/datasets/](https://huggingface.co/datasets/HiTZ/BertaQA)
503 [HiTZ/BertaQA](https://huggingface.co/datasets/HiTZ/BertaQA).

504 **Croissant Metadata.** The standard Croissant metadata of the dataset can be accessed at [https:](https://huggingface.co/api/datasets/HiTZ/BertaQA/croissant)
505 [//huggingface.co/api/datasets/HiTZ/BertaQA/croissant](https://huggingface.co/api/datasets/HiTZ/BertaQA/croissant).

506 **Additional Metadata.** To get dataset info: [https://datasets-server.huggingface.co/](https://datasets-server.huggingface.co/info?dataset=HiTZ/BertaQA)
507 [info?dataset=HiTZ/BertaQA](https://datasets-server.huggingface.co/info?dataset=HiTZ/BertaQA) To get the number of rows and size in bytes: [https://](https://datasets-server.huggingface.co/size?dataset=HiTZ/BertaQA)
508 datasets-server.huggingface.co/size?dataset=HiTZ/BertaQA. To get statistics for each
509 language check [https://datasets-server.huggingface.co/statistics?dataset=HiTZ/](https://datasets-server.huggingface.co/statistics?dataset=HiTZ/BertaQA&config=eu&split=test)
510 [BertaQA&config=eu&split=test](https://datasets-server.huggingface.co/statistics?dataset=HiTZ/BertaQA&config=eu&split=test) for Basque and [https://datasets-server.huggingface.](https://datasets-server.huggingface.co/statistics?dataset=HiTZ/BertaQA&config=en&split=test)
511 [co/statistics?dataset=HiTZ/BertaQA&config=en&split=test](https://datasets-server.huggingface.co/statistics?dataset=HiTZ/BertaQA&config=en&split=test) for English.

512 **Data hosting.** We host the dataset on HuggingFace at [https://huggingface.co/datasets/](https://huggingface.co/datasets/HiTZ/BertaQA)
513 [HiTZ/BertaQA](https://huggingface.co/datasets/HiTZ/BertaQA). The data can be viewed in the HuggingFace dataset viewer. It can be loaded with the
514 HF dataset loader. It can also be downloaded directly by cloning the HF repository or downloading
515 individual files. It can be loaded directly as it is in the standard JSONL format.

516 **Code hosting.** We host the code to reproduce experiments on GitHub: [https://github.com/](https://github.com/juletx/BertaQA)
517 [juletx/BertaQA](https://github.com/juletx/BertaQA). All instructions and code to reproduce our experiments are included in the repo.

518 **Licensing.** We do not own any of the text from which this data has been extracted. We license the
519 curation and translation of the dataset under CC-BY 4.0. The code is licensed under MIT license.

520 **Maintenance plan.** We will ensure that the dataset and code are available for a long time. We are
521 committed to maintaining the dataset and code to address any issues. We will actively monitor issues
522 in the HuggingFace and Github repositories.

523 **Author statement.** We bear all responsibility in case of violation of rights.