

Model Description

- Stable Beluga 2 is a Llama2 70B model finetuned on an Orca style Dataset. This repository contains the model from the stabilityai/StableBeluga2 repository with the following changes:
- **Storing weights in bfloat16 instead of float32.** This leads to 2x smaller files and a small quality loss, which is not significant compared to the loss caused by NF4 quantization used in Petals by default.
  - **Storing weights in small shards.** Each transformer block is stored in its own shard (1.71 GB each). The input and output embeddings and adjacent layernorms are in a separate shard (1.05 GB) too. This way, Petals clients and servers don't have to download any excess data besides the layers they actually use.
  - **Using Safetensors instead of Pickle.** This allows faster loading with smaller RAM requirements.

Model Details

Developed by: Stability AI  
Model type: Stable Beluga 2 is an auto-regressive language model fine-tuned on Llama2 70B.  
Language(s): English  
Library: HuggingFace Transformers  
License: Fine-tuned checkpoints (Stable Beluga 2) is licensed under the STABLE BELUGA NON-COMMERCIAL COMMUNITY LICENSE AGREEMENT  
Contact: For questions and comments about the model, please email lm@stability.ai

Training Dataset

Stable Beluga 2 is trained on our internal Orca-style dataset. Training data is a synthetic dataset that was created to enhance the small model's reasoning abilities. The dataset comprises a diverse collection of tasks aimed at training AI models across various domains, focusing on cautious reasoning and alignment with ethical guidelines. It includes approximately 602,000 zero-shot queries grouped into 23 categories and 126 sub-categories, each sharing a common instruction format to promote consistency. The dataset also features 55,000 few-shot samples to encourage the model's ability to learn from context, around 160,000 math problems sourced from a variety of existing datasets, and 2,000 synthetically generated conversations between doctors and patients designed to test the model's specialized skills.

Intended Uses

The pretrained-only model can be used for prompting for evaluation of downstream tasks as well as text generation. In addition, the model can be fine-tuned on a downstream task. For all other checkpoints, please have a look at the model hub.

Summary of Model Risks by Example Uses

We identified 14 potential model risks and 18 mitigation strategies (M) for 4 potential model uses					
		Recommending personalized content	Creating personalized ad campaigns	Personalizing learning experiences	Assisting in medical diagnoses
		(U1)	(U2)	(U3)	(U4)
Risks					
R1	Reflects offensive or biased content in model generated text despite finetuning	<div></div>	<div></div>	<div></div>	<div></div>
R2	Exhibits unreliable, unsafe or other undesirable behaviors	<div></div>	<div></div>	<div></div>	<div></div>
R3	Limits accuracy in areas underrepresented in the training dataset	<div></div>	<div></div>	<div></div>	<div></div>
R4	Facilitates the spread of misinformation by fabricating content making it unreliable for critical decisions	<div></div>	<div></div>	<div></div>	<div></div>
R5	Acts as 'black boxes' making it difficult to comprehend rationale behind outputs	<div></div>	<div></div>	<div></div>	<div></div>
R6	Carries biases present in the source data	<div></div>	<div></div>	<div></div>	<div></div>
R7	Exhibits limited real-world understanding	<div></div>	<div></div>	<div></div>	<div></div>
R8	Fails to cover all scenarios	<div></div>	<div></div>	<div></div>	<div></div>
R9	Harms economic interests of content creators by using their work without compensation	<div></div>	<div></div>	<div></div>	<div></div>
R10	Harms individuals' reputations by potentially revealing private medical information	<div></div>	<div></div>	<div></div>	<div></div>
R11	Increases psychological harm by exposing users to graphic and explicit content	<div></div>	<div></div>	<div></div>	<div></div>
R12	Infringes data protection laws by using sensitive data without proper authorization	<div></div>	<div></div>	<div></div>	<div></div>
R13	Undermines trust in AI systems by mishandling sensitive personal data	<div></div>	<div></div>	<div></div>	<div></div>
R14	Underperforms in non-English languages	<div></div>	<div></div>	<div></div>	<div></div>

Details of the Example Uses

- U1

Purpose: **Recommending personalized content**  
Capability: Analyzing preferences for suggestions  
AI User: Streaming platforms  
AI Subject: Content consumers  
Domain: Recommender Systems and Personalization
- U2

Purpose: **Creating personalized ad campaigns**  
Capability: Analyzing user behavior from social media posts  
AI User: Marketing agencies  
AI Subject: Consumers  
Domain: Marketing and Advertising
- U3

Purpose: **Personalizing learning experiences**  
Capability: Analyzing student performance and tailoring content  
AI User: Educational platforms  
AI Subject: Students  
Domain: Education and vocational training
- U4

Purpose: **Assisting in medical diagnoses**  
Capability: Analyzing patient data and suggesting conditions  
AI User: Healthcare professionals  
AI Subject: Patients  
Domain: Health and Healthcare

Mitigations for Risks

- R1

**Reflects offensive or biased content in model generated text despite finetuning**

Capability risk

Representation and toxicity harms

M

 Perform safety testing

M

 Tune model to specific applications

M

 Prevent harmful responses

M

 Leverage various content moderation services

M

 Avoid using models unsuitable for your application
- R2

**Exhibits unreliable, unsafe or other undesirable behaviors**

Capability risk

Information and safety harms

M

 Perform safety testing

M

 Prevent harmful responses

M

 Run your own suite of tests

M

 Avoid using models unsuitable for your application

M

 Exercise caution when using models in production systems

M

 Be mindful of potential issues in generated responses
- R3

**Limits accuracy in areas underrepresented in the training dataset**

Capability risk

Representation and toxicity harms

M

 Correct through evaluation and fine-tuning prior to deployment

M

 Requires detailed studies for better quantification of risks

M

 Need additional analysis to assess potential harm or bias
- R4

**Facilitates the spread of misinformation by fabricating content making it unreliable for critical decisions**

Human interaction risk

Misinformation harms

M

 Ensure outputs are not hallucinations

M

 Avoid treating model outputs as sources of truth

M

 Use the model responsibly

M

 Avoid using models for applications that may cause harm

M

 Perform more rigorous measurement, understanding and mitigations
- R5

**Acts as 'black boxes' making it difficult to comprehend rationale behind outputs**

Capability risk

Human autonomy and integrity harms

M

 Requires detailed studies for better quantification of risks

M

 Acknowledge the important role of research and open source community
- R6

**Carries biases present in the source data**

Capability risk

Representation and toxicity harms

M

 Perform safety testing

M

 Correct through evaluation and fine-tuning prior to deployment

M

 Leverage various content moderation services

M

 Perform more rigorous measurement, understanding and mitigations
- R7

**Exhibits limited real-world understanding**

Capability risk

Misinformation harms

M

 Correct through evaluation and fine-tuning prior to deployment

M

 Requires detailed studies for better quantification of risks

M

 Perform more rigorous measurement, understanding and mitigations
- R8

**Fails to cover all scenarios**

Capability risk

Information and safety harms

M

 Perform safety testing

M

 Correct through evaluation and fine-tuning prior to deployment

M

 Exercise reasonable caution when using models in production
- R9

**Harms economic interests of content creators by using their work without compensation**

Systemic risk

Socioeconomic and environmental harms

M

 Acknowledge the important role of research and open source community

M

 Follow better regulations and standards from government and technology leaders
- R10

**Harms individuals' reputations by potentially revealing private medical information**

Systemic risk

Information and safety harms

M

 Leverage various content moderation services

M

 Follow better regulations and standards from government and technology leaders
- R11

**Increases psychological harm by exposing users to graphic and explicit content**

Human interaction risk

Representation and toxicity harms

M

 Leverage various content moderation services

M

 Avoid using models in applications causing harm or distress

M

 Perform safety testing

M

 Follow better regulations and standards from government and technology leaders
- R12

**Infringes data protection laws by using sensitive data without proper authorization**

Human interaction risk

Information and safety harms

M

 Leverage various content moderation services

M

 Follow better regulations and standards from government and technology leaders
- R13

**Undermines trust in AI systems by mishandling sensitive personal data**

Human interaction risk

Information and safety harms

M

 Leverage various content moderation services

M

 Follow better regulations and standards from government and technology leaders

M

 Benefit from safety guardrails
- R14

**Underperforms in non-English languages**

Capability risk

Representation and toxicity harms

M

 Correct through evaluation and fine-tuning prior to deployment

M

 Requires detailed studies for better quantification of risks

M

 Need additional analysis to assess potential harm or bias