

# BertaQA: How Much Do Language Models Know About Local Culture?

Julen Etxaniz, Gorka Azkune, Aitor Soroa, Oier Lopez de Lacalle and Mikel Artetxe  
HiTZ Center, University of the Basque Country UPV/EHU



**GitHub**

<https://github.com/juletx/BertaQA>

# 1. Introduction

- LLMs have extensive knowledge about the world, but most evaluations have been limited to global or anglocentric subjects.
- How well do these models perform on topics relevant to other cultures, whose presence on the web is not that prominent?
- To address this gap, we introduce BertaQA, a multiple-choice trivia dataset

## 2. BertaQA

- First dataset with annotated questions related to the Basque and global cultures.
- 4756 parallel examples in Basque and English.
- Divided into 8 categories and 3 difficulties.
- Local subset with questions pertinent to the Basque culture
- Global subset with questions of broader interest.

	Local Questions	Global Questions
<b>Basque and Literature</b>	What does the “Karmel” magazine specialize in? a) Bertsolarism <b>b) Basque culture in the past and the present</b> c) The life of the Carmelites	In which of these novels does the sea not appear? <b>a) “The Adventures of Tom Sawyer”</b> b) “Moby Dick” c) “Treasure Island”
<b>Geography and History</b>	Where’s Atxondo? <b>a) In Biscay</b> b) In Gipuzkoa c) In Navarre	Who was imprisoned in 1964? <b>a) Nelson Mandela</b> b) Mumia Abu Jamal c) Charles Ghankay

### 3. Main results in English

- Open and commercial models much worse in local than global.
- Bigger difference for open models.
- Performance on local and global correlated.
- Scaling differences open vs closed.
- Easier to improve on global questions.
- But this subset starts saturating for the strongest models.
- Resulting in bigger improvements on the local subset.

Model	Variant	Local	Global	$\Delta$
Random	N/A	33.33	33.33	0.00
GPT	3.5 Turbo	55.08	82.40	27.32
	4	69.88	91.43	21.55
	4 Turbo	<b>72.17</b>	<b>91.68</b>	<b>19.51</b>
Claude 3	Haiku	58.71	84.16	25.45
	Sonnet	58.33	86.41	28.08
	Opus	<b>71.91</b>	<b>91.85</b>	<b>19.94</b>
Llama 2	7B	41.54	64.34	<b>22.80</b>
	13B	43.61	70.36	26.75
	70B	<b>49.15</b>	<b>77.68</b>	28.53
Llama 3	8B	50.38	76.63	26.25
	70B	<b>59.56</b>	<b>84.74</b>	<b>25.18</b>
Qwen 1.5	7B	42.51	71.45	<b>28.94</b>
	14B	44.67	75.92	31.25
	72B	<b>54.70</b>	<b>83.99</b>	29.29
Yi	6B	44.25	73.20	<b>28.95</b>
	9B	43.87	75.00	31.13
	34B	<b>54.06</b>	<b>83.61</b>	29.55
Mistral	7B	47.50	74.16	26.66
	47B	<b>57.40</b>	<b>82.78</b>	<b>25.38</b>
Gemma	7B	<b>45.69</b>	<b>76.42</b>	<b>30.73</b>
Average	N/A	53.25	79.91	26.66

## 4. Local knowledge transfer from Basque to English

- Local models trained with continued pretraining in Basque improve on local questions in English.
- Local models become worse on global questions.
- Bigger degradation and smaller improvement for the smallest model.
- Previous conclusions incomplete, challenges curse of multilinguality.

Model	Local	Global	$\Delta$
Llama 2 7B	41.54	<b>64.34</b>	22.80
+ <i>eu train</i>	<b>47.72</b>	53.26	<b>5.54</b>
Llama 2 13B	43.61	<b>70.36</b>	26.75
+ <i>eu train</i>	<b>56.60</b>	67.47	<b>10.87</b>
Llama 2 70B	49.15	<b>77.68</b>	28.53
+ <i>eu train</i>	<b>62.61</b>	73.62	<b>11.01</b>

## 5. Comparison of English and Basque

- Worse results in Basque for most models
- Local models better at answering local questions in Basque and global questions in English.
- Knowledge transfer is not perfect across languages.
- Local and global knowledge not transferred completely.

Model	Variant	Local	Global	$\Delta$
Random	N/A	33.33	33.33	0.00
GPT	3.5 Turbo	47.25 (-7.83)	66.22 (-16.18)	<b>18.97</b>
	4	62.94 (-6.94)	85.91 (-5.52)	22.97
	4 Turbo	<b>69.46</b> (-2.71)	<b>89.21</b> (-2.47)	19.75
Claude 3	Haiku	58.21 (-0.50)	79.85 (-4.31)	21.64
	Sonnet	56.13 (-2.20)	83.24 (-3.17)	27.11
	Opus	<b>71.32</b> (-0.59)	<b>90.89</b> (-0.96)	<b>19.57</b>
Llama 2	7B	34.90 (-6.64)	37.08 (-27.26)	<b>2.18</b>
	13B	34.09 (-9.52)	43.77 (-26.59)	9.68
	70B	<b>37.39</b> (-11.76)	<b>54.22</b> (-23.46)	16.83
Llama 2 + eu train	7B	49.45 (+1.73)	50.79 (-2.47)	<b>1.34</b>
	13B	60.24 (+3.64)	65.47 (-2.00)	5.23
	70B	<b>64.85</b> (+2.24)	<b>72.24</b> (-1.38)	7.39
Llama 3	8B	42.60 (-7.78)	63.09 (-13.54)	<b>20.49</b>
	70B	<b>57.40</b> (-2.16)	<b>82.15</b> (-2.59)	24.75
Qwen 1.5	7B	35.96 (-6.55)	46.15 (-25.30)	<b>10.19</b>
	14B	37.31 (-7.36)	53.39 (-22.53)	16.08
	72B	<b>42.77</b> (-11.93)	<b>63.25</b> (-20.74)	20.48
Yi	6B	37.94 (-10.32)	46.45 (-22.99)	<b>8.51</b>
	9B	38.20 (-13.79)	49.21 (-21.70)	11.01
	34B	<b>41.03</b> (-6.31)	<b>60.41</b> (-26.75)	19.38
Mistral	7B	37.18 (-5.67)	51.17 (-25.79)	<b>13.99</b>
	47B	<b>43.61</b> (-13.03)	<b>61.08</b> (-23.20)	17.47
Gemma	7B	<b>41.84</b> (-3.85)	<b>65.89</b> (-10.53)	<b>24.05</b>
<b>Average</b>	N/A	47.92 (-5.64)	63.53 (-14.41)	15.61

## 6. Translate-test and self-translate

- For Llama 2, translate-test improves results when compared to Basque. Self-translate does not provide big improvements.
- For local models, translate-test worsens results. Self-translate is better than translate-test, still does not reach Basque.
- For Gemma, translation improves global and harms local a bit.
- Overall, translation is better for global, does not work well on local.

Model	Size	Method	Local	Global
Llama 2	7B	Translate-test	<b>37.44</b> (+2.54)	<b>55.35</b> (+18.27)
		Self-translate	33.80 (-1.10)	38.71 (+1.63)
	13B	Translate-test	<b>37.69</b> (+3.60)	<b>62.50</b> (+18.73)
		Self-translate	34.81 (+0.72)	46.11 (+2.34)
	70B	Translate-test	<b>42.68</b> (+5.29)	<b>71.03</b> (+16.81)
		Self-translate	39.85 (+2.46)	55.23 (+1.01)
Llama 2 + eu train	7B	Translate-test	35.79 (-13.66)	44.27 (-6.52)
		Self-translate	<b>44.37</b> (-5.08)	<b>50.04</b> (-0.75)
	13B	Translate-test	41.79 (-18.45)	59.36 (-6.11)
		Self-translate	<b>56.13</b> (-4.11)	<b>65.55</b> (+0.08)
	70B	Translate-test	46.28 (-18.57)	65.47 (-6.77)
		Self-translate	<b>60.15</b> (-4.70)	<b>70.48</b> (-1.76)
Gemma	7B	Translate-test	<b>41.67</b> (-0.17)	<b>69.19</b> (+3.30)
		Self-translate	<b>41.67</b> (-0.17)	67.68 (+1.79)

# Thank you!



[julen.etxaniz@ehu.eus](mailto:julen.etxaniz@ehu.eus)



<https://julenetxaniz.eus/en>



[@juletxara](https://twitter.com/juletxara)