

RespondeoQA: a Benchmark for Bilingual Latin-English Question Answering

Marisa Hudspeth,¹ Patrick J. Burns,² Brendan O'Connor¹

¹Manning College of Information & Computer Sciences, UMass Amherst

²Institute for the Study of the Ancient World, New York University

Motivation

- Generative LLMs' performance on Latin is underexplored, despite being the 5th most prevalent language in recent large pretraining corpora
- Existing multilingual QA benchmarks exclude Latin entirely
- We introduce ~7,800 QA pairs from pedagogical sources spanning 1823–2025, with rich metadata across 10 content types

Dataset Sources & Statistics

Data Source	Year(s)	Type	Number of Question Types per Source			
			MC	1-W SA	Long Ans	Total
Exercises in Latin Prosody and Versification	1823	OCR scan	0	0	122	122
Latin Grammar and Junior Scholarship Papers	1832	OCR scan	0	675	350	1025
Certamen	1996–2009	Digital (Word)	317	4540	970	5827
National Latin Exam	2015, 2020, 2025	Digital (PDF)	855	0	0	855
			1172	5215	1442	7829

Dataset Curation

- OCR & Extraction** — use Gemini-1.5-pro to perform OCR on PDF scans + digital text
- Question-Answer Alignment** — combination of regex + GPT-4o to align questions to their answers
- Metadata Classification** — format, content, language, multihop, constrained translations
- Multipart Splitting** — Break multipart questions into standalone items

Question Page	Answer Key	Expanded questions																					
<p>JUNIOR SCHOLARSHIP PAPERS. 95</p> <p>95. Mention a Latin noun of which the gender is due to meaning in spite of form; and one of which the gender is due to form in spite of meaning.</p> <p>96. Translate and explain the subjunctives in—(1) <i>Fuisti quae in mensuris histeret.</i> (2) <i>Clautus ex quo hanc dicit.</i> (3) <i>Facundus fuerit; non mihi persuaderet.</i></p> <p>97. Give the gender, acc. sing., and meaning of—<i>Arctus, Boreas, Dialis, Cybele, Atlantis, Septentrio, Tibur, Thieris.</i> How else is <i>Thieris</i> spelt?</p> <p>98. By what rule of concord is <i>Animus et sententia pupuli in legibus posita</i> correct? What rules of concord does it break?</p> <p>99. Put into Latin—(1) "It is all over with the state." (2) "He has performed the work, which he undertook (regis) to do, to my satisfaction." (3) "They obtained their request for peace." (4) "Sixteen years after the fall of Carthage."</p>	<p>95. value, about 2%. Thus <i>centum sestertia</i> = 100,000 <i>sestertia</i>. In calculating sums above 100,000, the numeral <i>centum</i> is used with <i>gen.</i> and <i>omnes</i> <i>millia</i> understood; thus, <i>centa sestertia</i> = 1,000,000 <i>sestertia</i>.</p> <p>96. 1. <i>Optus</i>; <i>legis</i>. 2. <i>Arctus</i>, <i>abl.</i>; <i>Arctus</i>, <i>abl.</i>; <i>pupuli</i>, <i>gen.</i> or <i>abl.</i>; <i>in-stans</i>, <i>acc. s4</i> (3); <i>clautus</i>, <i>double acc.</i>; <i>clautus</i>, <i>abl.</i> Past parts of <i>creo</i> and <i>oro</i>. 3. (1) "Having languished for what they wore on their heads." Part of what she said (<i>oratio obliqua</i>).</p> <p>97. 1. <i>Arctus</i>, <i>gen. ind. of arctus</i>, <i>noun, acc.</i> <i>Arctus</i>, <i>noun, key</i>; <i>clautus</i>, <i>gen. ind. of clautus</i>, <i>noun, acc.</i> <i>clautus</i>, <i>noun, key</i>; <i>pupuli</i>, <i>gen. ind. of pupulus</i>, <i>noun, acc.</i> <i>pupulus</i>, <i>noun, key</i>; <i>Arctus</i>, <i>gen. ind. of arctus</i>, <i>noun, acc.</i> <i>Arctus</i>, <i>noun, key</i>; <i>clautus</i>, <i>gen. ind. of clautus</i>, <i>noun, acc.</i> <i>clautus</i>, <i>noun, key</i>; <i>pupuli</i>, <i>gen. ind. of pupulus</i>, <i>noun, acc.</i> <i>pupulus</i>, <i>noun, key</i>.</p> <p>98. (1) <i>Actum est de republica.</i> (2) <i>Opus, quod se facturum regis, e meo sententia fecit.</i> (3) <i>Pax impetrata.</i> (4) <i>Sexti decimo anno post Carthagem deletam.</i></p>	<table border="1"> <thead> <tr> <th>id</th> <th>question</th> <th>answer</th> </tr> </thead> <tbody> <tr> <td>95.1.1</td> <td>Is the gender of Latin noun "optio" determined by meaning or form?</td> <td>meaning</td> </tr> <tr> <td>95.7.1</td> <td>Give the perfect form of "molo."</td> <td>molui</td> </tr> <tr> <td>95.7.2</td> <td>Give the supine form of "molo."</td> <td>molitum</td> </tr> <tr> <td>95.7.29</td> <td>Give the supine form of "sano."</td> <td>sanatum</td> </tr> <tr> <td>95.8.1</td> <td>Put into Latin—"It is all over with the state."</td> <td>Actum est de republica.</td> </tr> <tr> <td>95.8.4</td> <td>Put into Latin—"Sixteen years after the fall of Carthage."</td> <td>Sexti decimo anno post Carthagem deletam.</td> </tr> </tbody> </table>	id	question	answer	95.1.1	Is the gender of Latin noun "optio" determined by meaning or form?	meaning	95.7.1	Give the perfect form of "molo."	molui	95.7.2	Give the supine form of "molo."	molitum	95.7.29	Give the supine form of "sano."	sanatum	95.8.1	Put into Latin—"It is all over with the state."	Actum est de republica.	95.8.4	Put into Latin—"Sixteen years after the fall of Carthage."	Sexti decimo anno post Carthagem deletam.
id	question	answer																					
95.1.1	Is the gender of Latin noun "optio" determined by meaning or form?	meaning																					
95.7.1	Give the perfect form of "molo."	molui																					
95.7.2	Give the supine form of "molo."	molitum																					
95.7.29	Give the supine form of "sano."	sanatum																					
95.8.1	Put into Latin—"It is all over with the state."	Actum est de republica.																					
95.8.4	Put into Latin—"Sixteen years after the fall of Carthage."	Sexti decimo anno post Carthagem deletam.																					

- Answer Refinement** — Simplify verbose answers for easier evaluation

Question Page	Answer Key	Simplified answers						
<p>EXERCISES. 2.</p> <p>Pāterjnā rūjā bōbūs exjēret sūis, Sōlūtūs omjni fājnōre.</p>	<p>CHAPTER I. FEET.</p> <p>The fifth foot in the first line, and the third in the second line, are spondees, all the other feet are iambs.</p>	<table border="1"> <thead> <tr> <th>id</th> <th>answer</th> </tr> </thead> <tbody> <tr> <td>app.1.2.1</td> <td>iambus, iambus, iambus, iambus, spondee, iambus</td> </tr> <tr> <td>app.1.2.2</td> <td>iambus, iambus, spondee, iambus</td> </tr> </tbody> </table>	id	answer	app.1.2.1	iambus, iambus, iambus, iambus, spondee, iambus	app.1.2.2	iambus, iambus, spondee, iambus
id	answer							
app.1.2.1	iambus, iambus, iambus, iambus, spondee, iambus							
app.1.2.2	iambus, iambus, spondee, iambus							

- Reference Translations** — create multiple explicit reference translations

"THERE WERE SIX (FEMALE) BEARS IN THE FOREST / SIX (FEMALE) BEARS WERE IN THE FOREST."	<ol style="list-style-type: none"> "six bears were in the forest" "there were six bears in the forest" "there were six female bears in the forest" "six female bears were in the forest"
--	--

- Filtering** — Remove unanswerable, diagram-dependent, and ambiguous questions

Question Taxonomy

10 content types × 4 language pair combinations. Top rows are knowledge-based; bottom rows are skill-based.

Content Type	Question Language → Answer Language				Total
	En → En	La → En	En → La	La → La	
Geography	73	1	107	1	182
History	253	0	685	3	941
Literature	85	2	311	0	398
Mythology	230	1	1283	8	1522
Vocabulary	698	9	733	21	1461
Grammar	192	122	894	106	1314
Lit. Devices	28	1	27	0	56
Read. Comp.	352	7	26	0	385
Scansion	21	20	59	42	142
Translation	854	91	483	0	1428
Total	2786	254	4610	179	7829

Key Findings

- All models struggle more with skill-based questions than knowledge-based ones
- LLaMA 3 leads overall (71.9%); reasoning models (QwQ, o3-mini) show advantages only on literary devices & scansion
- Best scansion accuracy is only 25% (o3-mini): significant room for improvement
- Question language matters: the reasoning models perform better on scansion questions asked in Latin, whereas LLaMa does better on English questions

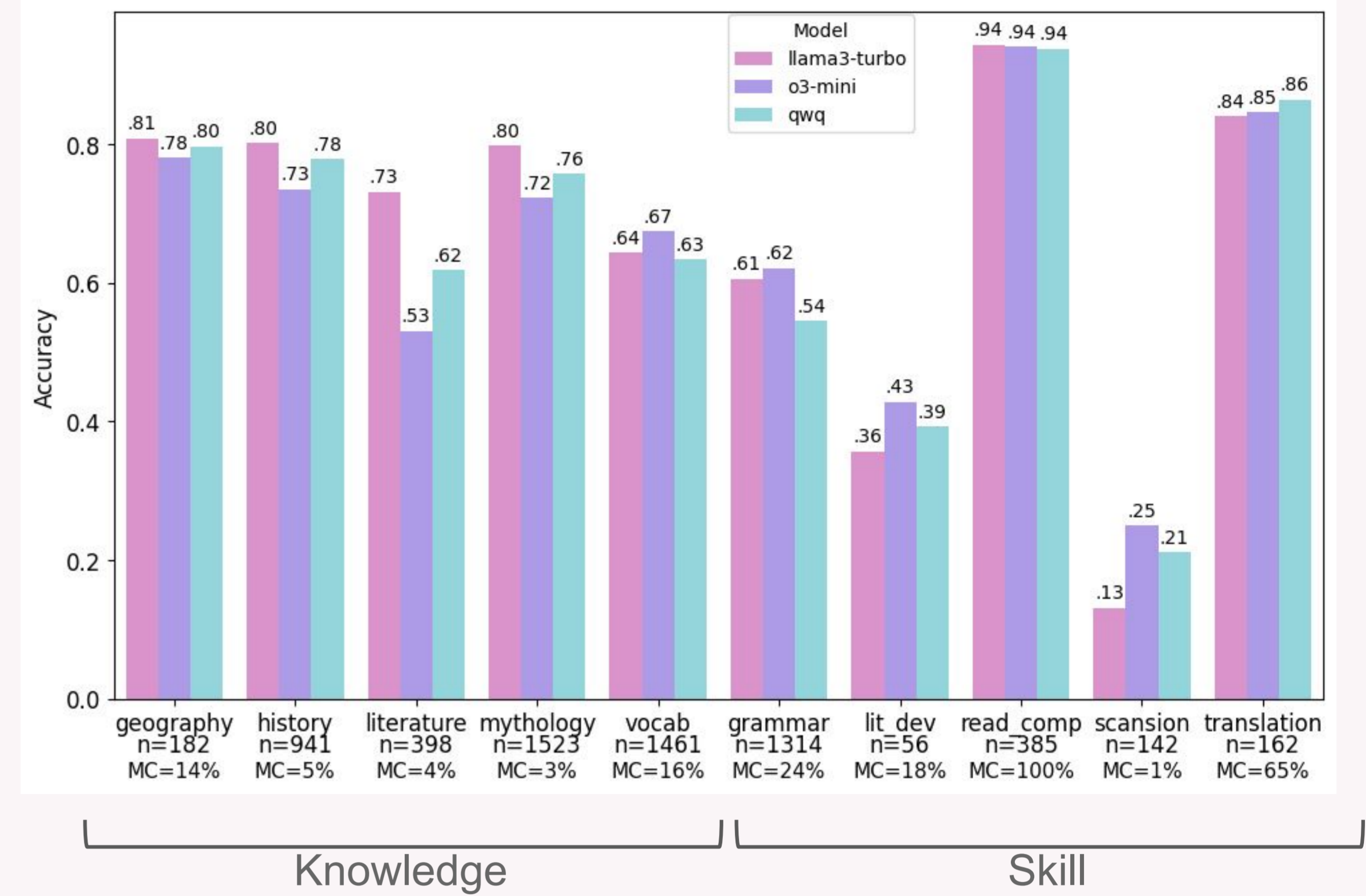
Results: Overall Accuracy

All models: knowledge > skill.

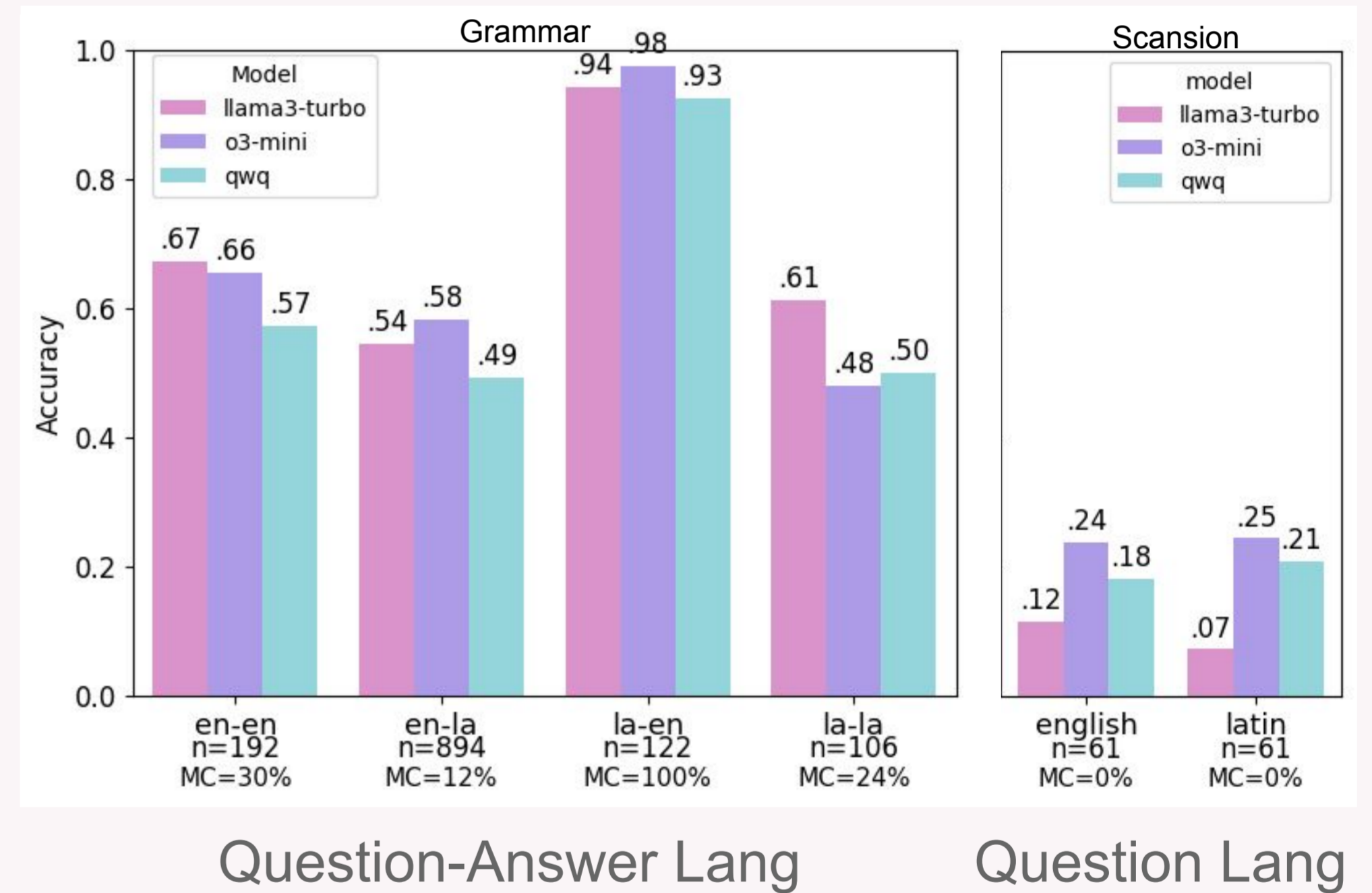
MC accuracy ~90% across all models.

	Overall	Knowledge	Skill
LLaMA 3	71.92	74.27	64.77
QwQ	68.11	71.08	61.63
o3-mini	68.61	69.46	66.76

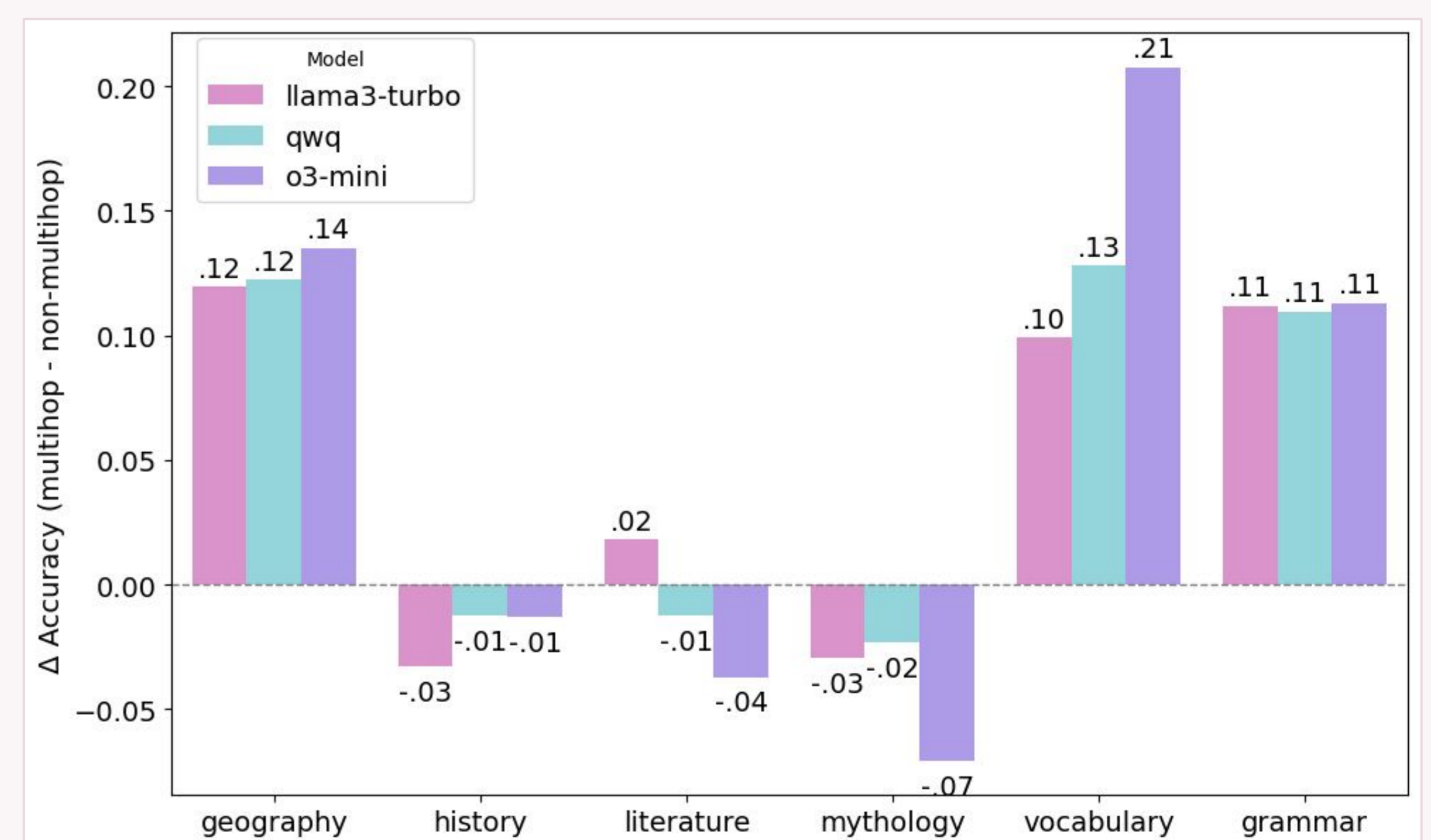
Accuracy by Content Type



Effect of Question Language



Multihop Question Accuracy



Translation Results (BLEU)

Setting	Llama 3		QwQ		o3-mini	
	Lat	Eng	Lat	Eng	Lat	Eng
Unconst.	25.50	45.41	18.53	39.91	23.42	43.29
Const.	20.88	37.46	21.52	17.95	27.25	34.68
Overall	23.71	45.25	19.53	39.27	24.67	43.14