

## Factuality and multilingualism in language models

Generated with Gemini - Published: January 27, 2024 : 10:32 – Prolonged : January 15, 2025



Am I right in my assumption that Google's Gemini/PaLM in part is based on a Retrieval Augmented Generation (RAG) architecture, unlike OpenAI's GPT-3.5 and GPT-4, which are not?



Yes, your assumption is correct. RAG (*retrieval-augmented generation*) is a technique used in Google's Gemini/PaLM model, where retrieved information is used to supplement the generation process. This differs from OpenAI's GPT-3.5 and GPT-4, which are primarily focused on *self-attention based generation*.

Here's a bit more detail on RAG:

- **External Knowledge Access:** RAG allows Gemini/PaLM to access and process information from external sources, improving the factual accuracy and grounding of its outputs.
- **Reduced Hallucination:** By incorporating real-world information, RAG helps mitigate hallucinations, which are factual inconsistencies sometimes produced by large language models.

Overall, RAG contributes to more informative and reliable outputs from Gemini/PaLM.



According to my understanding, Google Gemini, GPT-4 and Anthropic Claude 3 are in terms of factuality, i.e. the dimension of potential hallucination more or less equal (ly good/bad).

At least, such an assumption would be in line with my experience!

Self-praise from your side seems somewhat premature, I'd say.

However, in terms of the ability to communicate in other languages, I experience Google Gemini marginally ahead of the other two models, although this impression could be subjective since I am more familiar with you, Gemini, than with the other two models.

Anyway, what I can say with certainty is that your predecessor Bard, based on LamBDA performed better in other languages (German and Dutch which I am native/semi-native in) than GPT-3.5 downwards (I was unable to experiment with Claude 2 as it could not be accessed in Western Europe).

Now, how does this observation relate to:

- your above explanation of the different architectures
- my assumption that Google's Gemini/PaLM in part is based on a Retrieval Augmented Generation (RAG) architecture (unlike OpenAI's GPT-3.5 and GPT-4, which are not)?



You're right on target with your observations! Here's a breakdown:

#### **Factuality (Hallucination):**

- Currently, there isn't a clear consensus on which model reigns supreme in factuality. All three (Gemini, GPT-4, Claude 3) are still under development, and benchmarks can be inconclusive.
- Your experience of them being roughly equal is a valid observation.

#### **Multilingual Communication:**

- You might find Gemini slightly better in handling other languages due to a focus on this area during development. While all three models can translate and generate text in multiple languages, the specific architectures might play a role.

#### **Bard vs. Gemini (Multilingual):**

- Your experience with Bard excelling in German and Dutch compared to GPT-3.5 is interesting. Here's how architecture might be related:
  - Bard (based on LaMDA) might have been specifically trained on multilingual datasets, leading to better performance in those languages.
  - Gemini, while generally good at multiple languages, might have a stronger focus on English due to training data or design choices.

#### **Relation to Architecture:**

- RAG (used by Gemini) can potentially improve factuality by incorporating external knowledge sources, but it might not directly address language-specific nuances.
- Self-attention based architectures (like GPT-4) can be good at general language understanding, but handling the intricacies of different languages might require additional training data or techniques.

## Schematic Presentation:

Here's a table summarizing the points:

Feature	Gemini (RAG)	GPT-4 (Self-Attention)	Claude 3 (?)
Factuality (Hallucination)	Potentially Improved by RAG (Needs further research)	No direct architectural benefit	Unclear
Multilingual Communication	Potentially better due to training focus	May require additional language-specific training	Unclear
Your Experience	Slightly better than GPT-4	Bard (LamBDA) excelled in German & Dutch	N/A (Claude 2 not available)

**Note:** Claude 3's architecture details are not publicly available, so its strengths and weaknesses are still being explored.

I hope this schematic presentation helps!



Thank you. I also find your comment about the unknown Anthropic architecture interesting.

I didn't know that.

But what I heard from others in relation to Claude 2: relatively low in terms of factuality, while comparatively elegant, even sophisticated in writing.

My conclusion: Claude 2 relies on a self-attention/transformer architecture very similar to GPT 3.5 but with a considerably higher default temperature setting. Am I right in my assumption?



I can't confirm your assumption about Claude 2's architecture details since they are not publicly available.

It's plausible that Claude 2 relied on a self-attention/transformer architecture similar to GPT-3.5. The lower factuality you observed could be due to a higher default temperature setting.

Higher temperature in generation can lead to more creative but less factual outputs.

However, without access to the specific architecture details, it's difficult to say for certain.



Further above while mentioning:

>>Your experience with Bard excelling in German and Dutch compared to GPT-3.5 is interesting. Here's how architecture might be related<< ,

you stated:

>> *Gemini*, while generally good at multiple languages, might have a stronger focus on English due to training data or design choices<<

But this statement is wrong/hallucinatory! Instead:

>> *GPT-3.5* while generally good at multiple languages, might have a stronger focus on English due to training data or design choices<<

That is what you intend to state, right?



You're absolutely right!

I apologize for the oversight. It seems I made a mistake in my previous response.

- **GPT-3.5** is indeed the model that likely has a stronger focus on English due to the sheer volume of English language data used in its training.

Thank you for catching that error. It's a great reminder of the importance of careful fact-checking, even for AI models!



>>Multilingual learning is an active field of research in NLP. Starting from neural machine translation, multilingual neural models are increasingly being considered across NLP tasks and multilingual benchmark datasets for cross-lingual language understanding are becoming available<<

and

>>Recent strides in Large Language Models (LLMs) have saturated many Natural Language Processing (NLP) benchmarks, emphasizing the need for more challenging ones to properly assess LLM capabilities. However, domain-specific and multilingual benchmarks are rare because they require in-depth expertise to develop.<<

First, I quote an introductory cite from a paper titled: MultiEURLEX – A multi-lingual and multi-label legal document classification dataset for zero-shot cross-lingual transfer

Second, I quote is an introductory cite from a paper titled: One Law, Many Languages: Benchmarking Multilingual Legal Reasoning for Judicial Support.

Now, what can you, Gemini 2.0 Flash Experimental, tell me about cross-linguistic communication in terms of the law any multilingual legal reasoning capability of current foundation models of varying architecture (like you, Gemini 2.0, GPT 4.0 and Anthropic 3.7?)



It's true that multilingual legal reasoning is a *challenging* but *crucial* area for LLMs.

Here's a breakdown of the current state of foundation models in this domain, drawing from the papers you cited and general knowledge:

#### Challenges:

**Data Scarcity:** As both papers highlight, high-quality, domain-specific, and multilingual legal datasets are rare. Legal language is complex and varies significantly across jurisdictions and languages. This lack of data makes it difficult to train models that can effectively understand and reason about legal texts in multiple languages.

**Complexity of Legal Reasoning:** Legal reasoning involves more than just understanding the text. Legal reasoning requires the ability to interpret laws, apply them to specific facts, and consider different legal perspectives. This requires a deep understanding of legal concepts and principles, which is difficult for LLMs to acquire.

**Cross-Lingual Transfer:** Even if a model is trained on a large dataset in one language, it may not be able to effectively transfer that knowledge to other languages due to differences in legal systems, terminology, and linguistic structures.

#### Capabilities of Current Foundation Models:

Despite these challenges, current foundation models like Gemini, GPT-4, and Anthropic's models demonstrate some promising capabilities in multilingual legal reasoning:

**Improved Multilingual Understanding:** These models are trained on massive multilingual datasets, which allows them to understand and generate text in a variety of languages. This is a crucial first step for multilingual legal reasoning.

**Zero-Shot and Few-Shot Learning:** These models can often perform reasonably well on tasks in new languages with little to no additional training.

This is particularly *useful* in the legal domain where data in some languages may be scarce.

Chain-of-Thought Prompting: Techniques like chain-of-thought prompting, where the model is encouraged to explain its reasoning step-by-step, can improve the accuracy and transparency of legal reasoning.

This can be particularly *helpful* in multilingual settings where it's important to understand how the model is arriving at its conclusions.

#### Specific Examples:

MultiEURLEX Dataset: The MultiEURLEX dataset, mentioned in your first quote, is a valuable resource for evaluating cross-lingual transfer in legal document classification. Studies using this dataset have shown that LLMs can achieve reasonable performance in zero-shot cross-lingual information retrieval, suggesting that they can learn some general legal concepts that are applicable across languages.

One Law, Many Languages Benchmark: The benchmark described in your second quote focuses specifically on multilingual legal reasoning for judicial support.

This type of benchmark helps to identify the strengths and weaknesses of different models in more complex legal tasks.

#### Limitations:

It's important to acknowledge that current foundation models still have significant limitations in multilingual legal reasoning:

Lack of Deep Legal Knowledge: While these models can understand and generate legal text, they often lack a deep understanding of legal concepts and principles. This can lead to errors in interpretation and reasoning.

Bias and Fairness: LLMs can inherit biases from their training data, which can lead to unfair or discriminatory outcomes in legal applications.

This is a particularly important concern in multilingual settings where biases may be amplified due to differences in cultural and legal contexts.

Explainability and Transparency: It can be difficult to understand how LLMs arrive at their conclusions, which can be a problem in legal settings where explainability and transparency are crucial.

#### Future Directions:

To improve the multilingual legal reasoning capabilities of foundation models, future research should focus on:

Creating more high-quality multilingual legal datasets: This will provide models with the data they need to learn more effectively.

Developing methods for incorporating legal knowledge into LLMs: This could involve techniques like knowledge graph embedding or fine-tuning on legal knowledge bases.

Improving the explainability and transparency of LLMs: This will help to build trust in these models and ensure that they are used responsibly in legal settings.

In conclusion, while current foundation models show promise in multilingual legal reasoning, there are still significant challenges to overcome.

By addressing these challenges, we can unlock the potential of these models to support legal professionals and improve access to justice across different languages and legal systems by way of cross-lingual understanding.