**GOOGLE GEMINI**

# Factuality and multilingualism of different language models

Generated with Gemini - Published: January 27, 2024 10:32

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.

If I understand it correctly, 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, as far as the ability to communicate in other languages is concerned, I find Google Gemini slightly, but only slightly, 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.

What I can say with certainty, however, is that your predecessor Bard, based on LamBDA presented 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 which was not available Western Europe).

Now, how does this observation relate to your above explanation of the different architectures?

Would you be able to present this answer to me in a schematic form in addition to text form?

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 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.

One last thing, further about 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.

>>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!