

## **Understanding AI models: A comparative overview**

Before we start coding, let's dive into the different AI models we'll be working with. Understanding their strengths, weaknesses, and use cases is crucial for building effective GenAI applications.

### **Llama 3**

Llama 3 is the latest iteration in the Llama series, building upon the success of Llama 2.

#### **Strengths:**

- Improved performance over Llama 2 in many tasks
- Enhanced context understanding and generation capabilities
- Better handling of nuanced prompts

#### **Weaknesses:**

- Higher computational requirements compared to Llama 2
- Often requires more fine-tuning or prompt engineering compared to other models.

#### **Best use cases:**

- Advanced language understanding and generation tasks
- Applications requiring up-to-date knowledge and improved reasoning

### **Granite**

Granite is IBM's advanced language model, part of the watsonx.ai platform.

#### **Strengths:**

- Optimized for enterprise use cases
- Strong performance in business and technical domains
- Integration with IBM's ecosystem of tools and services

#### **Weaknesses:**

- May be less versatile for general-purpose tasks compared to open-source alternatives
- Access and usage may be more restricted compared to open-source models

#### **Best use cases:**

- Enterprise-level applications
- Specialized business and technical tasks
- Integration with other IBM services

## **Mixtral**

Mixtral by Mistral AI uses a Mixture of Experts (MoE) setup, where each layer has 8 specialized "experts," selecting the best ones for each task.

### **Strengths:**

- Efficient as it activates only the necessary "experts," making it resource-efficient for diverse tasks.
- High adaptability due to specialized experts, allowing fine-tuning for specific needs.
- Performs well on both general and specialized tasks without increasing computational costs drastically.

### **Weaknesses:**

- The MoE structure can add complexity in deployment and model interpretation.
- Over-specialization in MoE systems can lead to overfitting, where experts trained on narrow data subsets perform poorly on new data, lowering overall system accuracy.
- As a newer architecture, it may have fewer pre-trained variants and community resources.

### **Best use cases**

- Adaptive systems applications needing flexibility to handle various task types with optimized resource usage.
- Customizable tasks: Scenarios where fine-tuning for domain-specific tasks is critical, such as specialized industry applications.

### **Performance considerations**

When choosing between these models, consider:

1. **Speed:** Smaller models like Llama 3.2 1B might be faster but less capable. Larger models like Llama 3 70B or Granite might be slower but more powerful.
2. **Accuracy:** Generally, larger models tend to be more accurate, but this can vary depending on the specific task.

3. **Cost:** Larger models and proprietary models like OpenAI's GPT may incur higher usage costs.
4. **Latency:** Consider the response time requirements of your application. Smaller models or edge-deployable versions might be preferable for low-latency applications.
5. **Specialization:** Some models might perform better for specific domains or tasks. Granite, for instance, might excel in business-oriented tasks.

Understanding these trade-offs will help you make informed decisions when implementing AI in your applications.