

Dr. Ina Mittelstädt

(Stabsstelle Organisations- und Personalentwicklung)

# **ChatGPT, Qwen, Apertus – Navigating the Jungle of AI Models (Chatbots)**

KI-Werkstattgespräch zur Mittagspause



# Where is ChatGPT?

- Apertus 70B Instruct 2509
- DeepSeek R1 Distill Llama 70B
- Devstral 2 123B Instruct 2512
- Gemma 3 27B Instruct
- GLM-4.7
- InternVL 3.5 30B A3B
- Llama 3.1 SauerkrautLM 70B Instruct
- MedGemma 27B Instruct
- Meta Llama 3.1 8B Instruct

„Models“ =  
Large Language Models  
(LLMs)

Currently: around 110 LLMs  
([wikipedia](#))

# Which one is the right one for me?



## General differences: Proprietary or open-weight/open-source?

### Proprietary

- ChatGPT
- Claude
- CoPilot
- Gemini
- NotebookLM

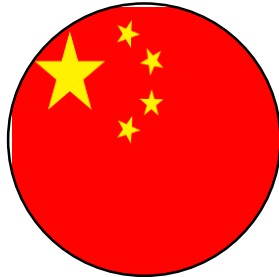
### Open Weights/ Open Source

- Apertus
- Teuken
- Mistral
- OLMo
- GPT OSS 120B
- Llama
- Gemma

# General differences: Developers' backgrounds



OpenAI (ChatGPT)  
Anthropic (Claude)  
Google (Gemini)  
Microsoft (CoPilot)  
Meta (Llama) ...



Alibaba (Qwen)  
Deepseek  
OpenGVLab (InternVL)  
Moonshot (Kimi)...



Mistral (Mistral, Devstral)  
Swiss AI (Apertus)  
OpenGPT-X (Teuken)...

# General Differences: Specialization

## **Instruct:**

- Quick answers, retrieving knowledge, doing routine tasks

## **Reasoning/Thinking:**

- Spending more time on the tasks, breaking them down into individual steps, weighing the options (good for solving problems)

## **Multimodal:**

- Can process inputs in different modes simultaneously (image + text, image + audio, etc.)

## **Arcana (RAG):**

- Operates based on specific provided knowledge resources (documents, databases, etc.)

# How do I choose the right model?

- Data Protection



# How do I choose the right model?

## Possible reasoning: Data Protection

### Terms of Use ChatGPT:

*Our use of content.* We may use Content to provide, maintain, develop, and improve our Services, comply with applicable law, enforce our terms and policies, and keep our Services safe.

### EU Privacy Policy ChatGPT:

*Usage Data:* We collect information about your use and activity across the Services, such as the types of content that you view or engage with, the features you use and the actions you take, when you submit feedback to a model response, the people with whom you interact.

... We may share your Personal Data ... with government authorities, industry peers, or other third parties in compliance with the law ...

# How do I choose the right model?

## **Possible reasoning:** Data Protection

When using AI in a professional context (while working for the EUF), do NOT enter personal data, confidential information, or legally protected documents into commercial AI tools (like ChatGPT, Claude, Perplexity)! (Except university leadership states otherwise at some point).

When it comes to data protection, Academic Cloud's ChatAI is a safe choice!



# How do I choose the right model?

- Data Protection
- Intended use

# How do I choose the right model?

**Possible reasoning:** Intended use

## Instruct

- Retrieve information
- Write texts (emails, academic papers, meeting minutes, etc.)
- Explanations
- Summarize documents
- Create questions
- Overviews

## Thinking/Reasoning

- Strategic planning
- Complex problems
- Generating and discussing challenging ideas (e.g., research)
- Developing an argument

## Special Applications

- Image Processing
- Mathematics
- Coding
- Tool Use
- Projects
- Medicine

# How do I choose the right model?

- Data Protection
- Intended use
- Performance

# How do I choose the right model?

**Possible reasoning:** Performance

[Available Models :: Documentation for HPC](#)

<https://llm-stats.com/benchmarks>

<https://artificialanalysis.ai/>



# How do I choose the right model?

- Data Protection
- Intended use
- Performance
- Number of parameters

# How do I choose the right model?

**Possible reasoning:** Number of parameters

For instance: Apertus **70B** Instruct, Gemma 3 **27B** Instruct, Meta Llama 3.1 **8B**

Or: 8B, 27B, 30B, 70B, 123B, 397B, 675B... (B = Billion)

What are parameters?

- NOT: Stored or compressed facts (hence: more parameters does NOT mean more knowledge)
- RATHER: The set of learned semantic relationships
- More parameters = greater accuracy in complex tasks (but also: higher energy consumption)
- For many tasks, just a few parameters are sufficient!



# How do I choose the right model?

- Data Protection
- Intended use
- Performance
- Number of parameters
- Context window

# How do I choose the right model?

**Possible reasoning:** Size of the Context window

Context window = How much information can the LLM keep in memory to continue working with it?

Exceeding it = "Forgetting" the previous information

Scale = Token (words or parts of words)

K = 1000

Information about models in ChatAI (in english):

[Available Models :: Documentation for HPC](#)



# How do I choose the right model?

- Data Protection
- Intended use
- Performance
- Number of parameters
- Context window
- Language

# How do I choose the right model?

**Possible reasoning:** Language (other than English)

Check <https://artificialanalysis.ai/models/multilingual> for benchmark results for several languages

Apertus has the mission to support small European languages!  
(f.i. trained in Swiss German and Romansh)



# How do I choose the right model?

- Data Protection
- Intended use
- Performance
- Number of parameters
- Context window
- Language
- Energy consumption

# How do I choose the right model?

**Possible reasoning:** Energy consumption

Smaller models consume less energy when generating responses

Instruct uses less energy than Reasoning/Thinking

☞ „A3B“ / „A10B“

large models that activate (use) only a small portion of the parameters (= a newer technical approach)

→ Saves energy

# How do I choose the right model?

- Data Protection
- Intended use
- Performance
- Number of parameters
- Context window
- Language
- Energy Consumption
- Personal taste/needs

## Conclusion

- Don't look for the perfect model for everything: All models have their strengths and weaknesses
- Try out different models – with different tasks
- New models are constantly being developed—don't get too used to just one!
- Set your own standards (and follow your own taste)
- ALL models make mistakes—never trust them blindly!

