# LLM PITFALLS

Severity  LOW ▼ ● ▲ HIGH

| Category | Pitfall | Severity | Mitigation | GPT 3.5 (0.0%) | GPT 4 (31.0%) | GPT o4-mini (58.6%) | Claude 3.5 (31.0%) | Claude 4 (58.6%) | Sonar (37.9%) |
|---|---|---|---|---|---|---|---|---|---|
| **Design** | Using generative AI as a first resort | ▲ | Consider more explainable, scrutable tech | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | No baseline or state-of-the-art performance | ● | Measure SOTA performance | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | No success criterion | ● | Ask user/customer what success means | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| **Instruction following** | Refusal to answer on unreasonable grounds | ▼ | Add context, remove ambiguity | Fail | Pass | Pass | Fail | Pass | Pass |
| | Misinterpretation, failure to clarify | ▲ | Instruct to clarify, add context, be explicit | Fail | Pass | Pass | Fail | Pass | Pass |
| | Failure to reject false or inconsistent query | ● | Avoid incorrect or inconsistent prompts | Fail | Fail | Fail | Pass | Pass | Pass |
| | Arbitrary execution of instructions in data | ▲ | Explicitly separate instructions and data | Fail | Fail | Fail | Fail | Pass | Fail |
| | Difficulty manipulating sub-token elements | ● | Pre-process sub-token elements | Fail | Fail | Pass | Fail | Fail | Fail |
| | Incorrectly formatted output | ● | Provide format schema or template | Fail | Fail | Fail | Fail | Pass | Fail |
| **Ambiguity & coherence** | Ambiguity in the response | ▲ | Prompt engineering, pre-check responses | Fail | Fail | Pass | Fail | Fail | Fail |
| | Incoherent rambling or glitching | ▼ | Reword prompt and retry | Fail | Fail | Pass | Not evaluated | Pass | Fail |
| | Inconsistent reasoning across conversation | ● | Shorten conversations, provide recaps | Fail | Pass | Pass | Not evaluated | Pass | Pass |
| | Irrelevance | ▼ | Prompt engineering, pre-check responses | Fail | Pass | Pass | Pass | Pass | Fail |
| **Response quality** | Simplistic (correct but poor) responses | ● | Role allocation, few-shot prompts | Fail | Pass | Pass | Pass | Fail | Pass |
| | Incompleteness | ● | Role allocation, few-shot prompts | Fail | Fail | Pass | Fail | Fail | Fail |
| | Faulty reasoning | ▲ | Chain-of-thought, tools, few-shot prompts | Fail | Pass | Pass | Pass | Pass | Fail |
| | Faulty premises | ▲ | Explicit chain-of-thought, break steps down | Fail | Fail | Fail | Fail | Fail | Fail |
| | Prompt sensitivity | ● | Add more context and more instruction | Fail | Fail | Pass | Pass | Pass | Fail |
| | Overthinking | ● | Provide more context about problem domain | Fail | Fail | Pass | Fail | Fail | Fail |
| | Overfit to similar task | ▼ | Emphasize differences from overfitted case | Fail | Fail | Fail | Fail | Pass | Fail |
| | Distraction | ▼ | Use neutral wording | Fail | Fail | Fail | Fail | Fail | Fail |
| | Verbose | ▼ | Specify minimal requirements | Fail | Fail | Fail | Fail | Fail | Fail |
| **Factual errors** | Incorrect factual recall | ▲ | RAG, knowledge graphs, pre-check response | Fail | Fail | Pass | Fail | Pass | Fail |
| | Fabricated or non-useful references | ▲ | RAG, knowledge graphs | Fail | Fail | Fail | Fail | Fail | Pass |
| | Coherent fiction | ▲ | RAG, knowledge graphs, pre-check response | Fail | Fail | Fail | Pass | Pass | Pass |
| | Out-of-date responses | ● | RAG, tools | Fail | Pass | Fail | Fail | Fail | Pass |
| | Overconfidence in certainty | ● | Disallow expressions of uncertainty | Fail | Fail | Pass | Fail | Fail | Fail |
| **Bias & ethics** | Explicit bias | ▲ | Improve alignment and system prompt | Fail | Fail | Pass | Pass | Pass | Pass |
| | Implicit bias | ▲ | Improve alignment, avoid leading signals | Fail | Fail | Fail | Fail | Pass | Fail |
| | Toxic or offensive responses | ▲ | Pre-check response | Fail | Fail | Fail | Fail | Pass | Fail |
| | Copyright infringement | ▲ | Knowledge graphs | Fail | Pass | Pass | Pass | Pass | Pass |
| | Cultural or ideological bias | ● | Fine-tuning, alignment | Fail | Pass | Pass | Pass | Fail | Pass |
| **Application** | No transparency about role of AI | ▲ | Be clear | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | Not collecting signals from users | ▼ | Collect signals | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | Not providing references | ● | Add references | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| **Implementation** | No rigorous evaluation protocol | ▲ | Create testing and reporting pipeline | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | Low AI literacy among users | ● | Education and culture | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | Overreliance on model responses | ▲ | Education, avoiding risky applications | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | Overreliance on human-in-the-loop | ● | Consultation, training, improved interface | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | No user training | ● | Training | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |
| | No consideration of security or ethics | ▲ | Proper governance and oversight | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated | Not evaluated |

A highly opinionated and non-exhaustive list of pitfalls aimed at developers of applications containing large language models,. There are more than 30 concerns here, but some of them overlap, and probably still others are missing entirely. This list is my own opinion, which is not necessarily shared by others at Equinor or anywhere else. It is a work in progress, you input will be welcomed and credited, please get in touch!

Abbreviations: RAG means 'retrieval augmented generation'.