

SplunkGPT

Unlocking LLM Superpowers

Andrew Gomez

Jake Coyne



whoami



Andrew Gomez

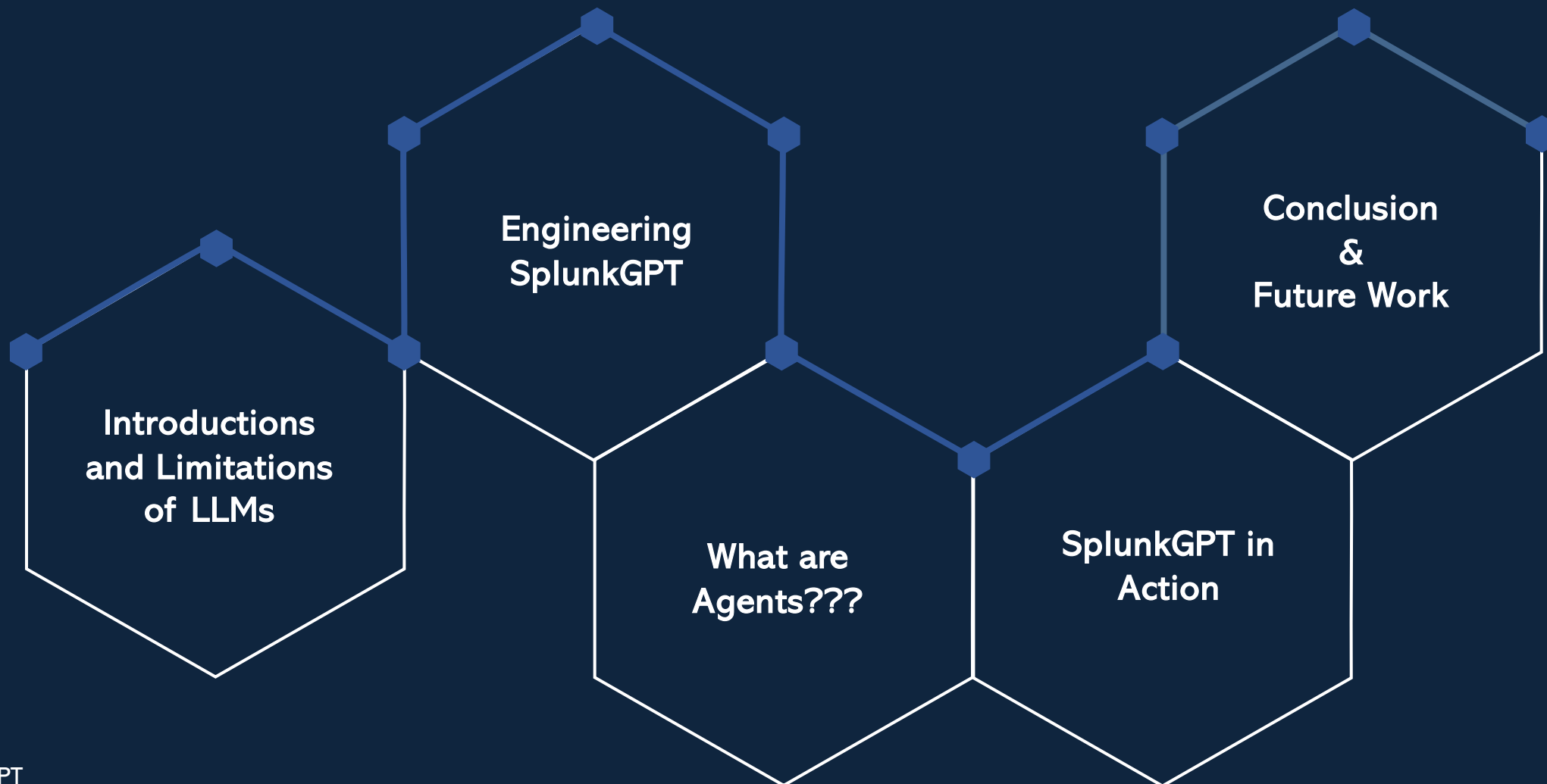
Offensive Operator
Sixgen



Jake Coyne

Offensive Operator
Sixgen

Agenda



Introduction and Limitations of LLMs



What are LLMs?

- Type of AI to predict new content
- Examples: Llama, Flan, Bart, GPT, etc...



Token Limits

- GPT (3) text prompt 4096 characters limit
- 4096 tokens per conversation
- Varies per LLM and model



Computational Costs

- 4K Model | \$0.0015 / 1K tokens | \$0.002 / 1K tokens
- 1 Token \approx 4 characters
- Quadratic increase in cost
- Tuning is \$\$\$
- Cost can also be energy if local LLM is used



Hallucinations

- Limited to a snapshot in time
- Black Box

A potential solution... Create an Agent!

Engineering SplunkGPT



Prompt Engineering Methods

- ZeroShot
- FewShot
- Chain-of-Thought
- Plan and Solve
- Reasoning and Acting



Memory Management

- Short Term Memory vs Long Term Memory



Vector Databases

- What is it?
- Why is it important?
- Algorithms



Tools for Interaction

- LangChain Integration with...
 - Google Search
 - Hugging Face
 - Dall-E



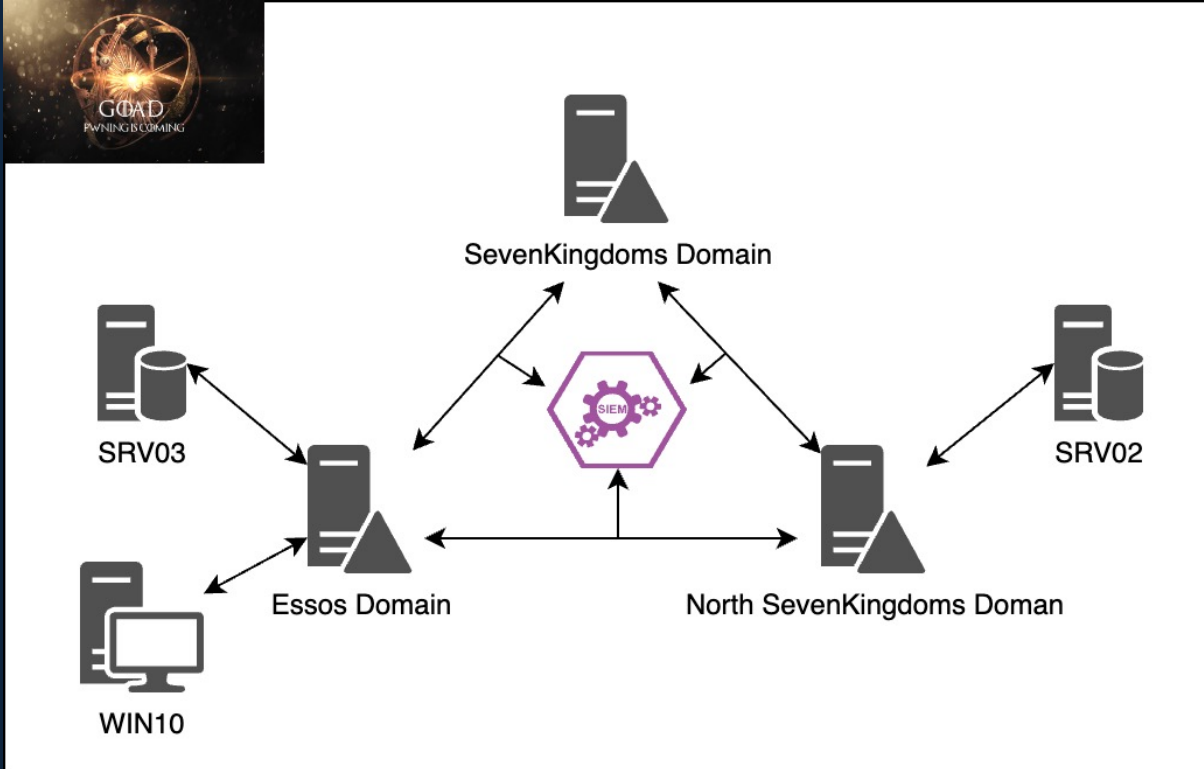
Function Calls

- Tools as functions
- Challenges Faced



Chaining it together as an Agent

- What is an agent?
- How is LangChains tying it all together?



DEMO TIME

References and Resources



References

- Greg Kamradt (Data Indy)
- BabyAGI
- AgentGPT
- SplunkAI
- OpenAI Documentation
- LangChain Documentation



Future Work

- Release Tool: <https://github.com/andrew-gomez/SplunkGPT>
- Refine Streamlit Chat Application
- Migration to a local LLM
- Reduce token size of prompts
- Continue testing and refining