



# SquirroGPT

10 Considerations for constructing a RAG – Checklist

---

Uncover Insights with AI-Apps from Squirro



---

# Agenda

10 Considerations for constructing a RAG

Retrieval Augmented Generation

Meet Squirro

---

# 10 Considerations for constructing a RAG

# 10 Considerations for constructing a RAG

#	Consideration	Check	Notes
1	<b>Data Access &amp; Life Cycle Management</b>		<ul style="list-style-type: none"> <li>• Manage the entire lifecycle of the information</li> <li>• Data is efficiently processed, enriched, available and eventually deleted</li> </ul>
2	<b>Data Indexing &amp; Hybrid Search</b>		<ul style="list-style-type: none"> <li>• Setup of hybrid search</li> <li>• Maintaining such a large-scale index over time</li> </ul>
3	<b>Enterprise Security &amp; Access Control at Scale</b>		<ul style="list-style-type: none"> <li>• Respecting complex Access Control Lists (ACL),</li> <li>• Secure setup of a RAG system</li> </ul>
4	<b>Chat User Interface</b>		<ul style="list-style-type: none"> <li>• Adaptable chat interface</li> <li>• Extendible chat interface for new task classes</li> </ul>
5	<b>Comprehensive System Interaction</b>		<ul style="list-style-type: none"> <li>• Management of interactions between Information retrieval, user interface, LLM and entailment check</li> <li>• Comprehensive Information Retrieval (IR) Stack</li> </ul>
6	<b>Prompt Engineering</b>		<ul style="list-style-type: none"> <li>• Creating an effective prompt service</li> <li>• Integrating adaptive mechanisms can help refine prompts based on real-time interactions and feedback</li> </ul>
7	<b>Chain of Reasoning</b>		<ul style="list-style-type: none"> <li>• Sophisticated user journey progression support</li> </ul>
8	<b>Enterprise Integration</b>		<ul style="list-style-type: none"> <li>• Integrating a RAG into an existing enterprise setup (SDK, API , etc.)</li> </ul>
9	<b>Continuous Operation</b>		<ul style="list-style-type: none"> <li>• Service Level requirements, organizational setup, etc.</li> <li>• Talent availability to operate such systems over time</li> </ul>
10	<b>Cost Considerations</b>		<ul style="list-style-type: none"> <li>• Find the good balance between the LLM and the IR stack</li> <li>• Total cost of ownership considerations</li> </ul>

— SQUIRRO GPT

# Retrieval Augmented Generation

*“Autoregressive LLMs are doomed*

*There will be better systems that are factual,  
non-toxic, and controllable*

*Marrying them with tools such as search  
engines is highly non-trivial”*

Yann Andre LeCun

Chief AI Scientist at Facebook & Silver Professor at  
the Courant Institute, New York University.

# Combine an Insight Cloud with Generative AI<sup>1</sup>

Generative AI models hallucinate  
**Answers come with the "why"**

Works with public data; only limited with company data  
**Answers are based on your own data**

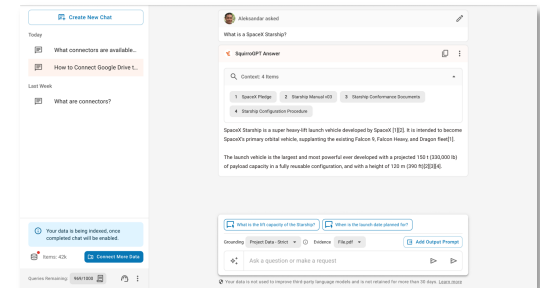
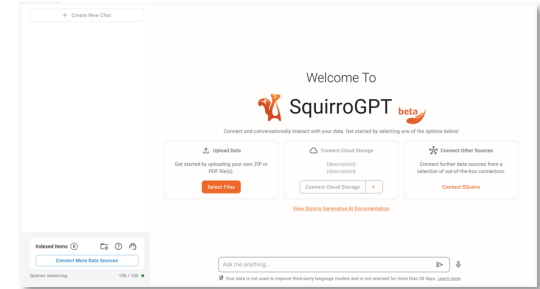
They have a disdain for enterprise security  
**Robust, secure, and respecting ACLs**

& more<sup>1</sup>



# An Enterprise ready Application Stack

- No Hallucination ✓ Evidence based:  
Answers come with the "why"
- Local content ✓ Your Company Data at Scale:  
Local context and trained LLMs
- Enterprise security ✓ Robust, secure, and performant  
Enterprise ACL setup, and more





# Retrieval Augmented Generation (RAG)

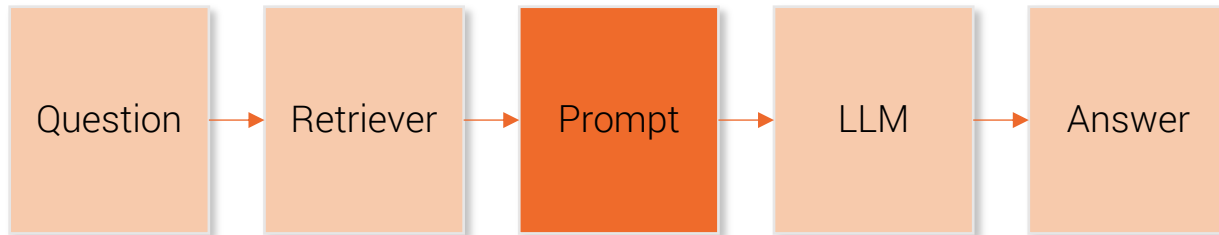
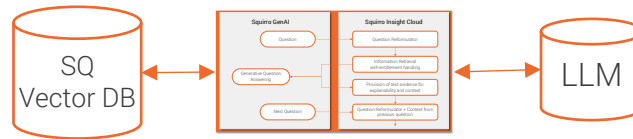
Load related context / information into "working memory" LLM context window



Google / Bing / Squirro  
Classic Enterprise Search  
(Retrieval)

Retrieve relevant information;  
Expose to LLMs of your choice

OpenAI, Cohere, Anthropic  
Large Language Model  
(Generation)



# Generative AI stack

## Application

End-to-End applications including models or Apps without proprietary models (ability to use a best of breed approach)

## Models

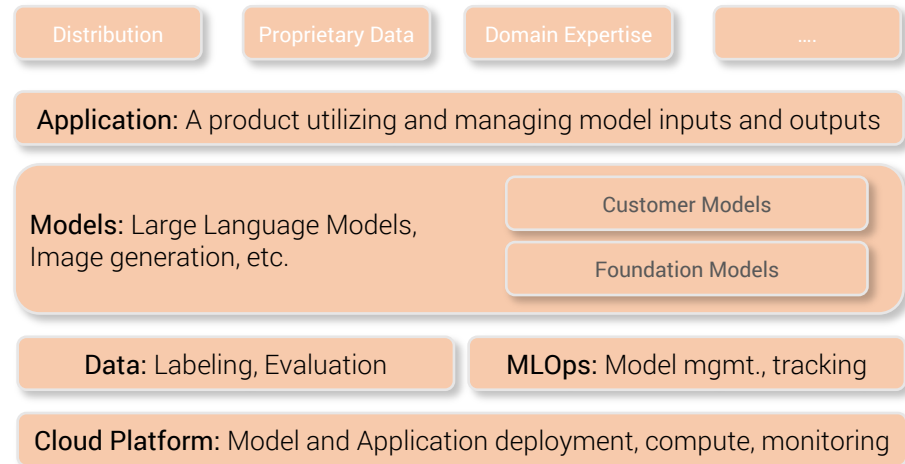
Splitting the model layer to differentiate between proprietary and open-source models as well between foundation and custom models

## Data & MLOps

Models derive their value from the data they're trained on. So, there's a need to factor in data and machine learning operations (MLOps)

## Cloud Platform

To run GPU based platform locally is hard; Time to switch to cloud platforms



---

# Meet Squirro

# Meet Squirro

## Global Presence

- Augmented Intelligence Software Provider
- Presence in Europe, United States and Asia
- Global presence and partner network

## A Different Approach to AI

- A powerful Insights Engine for Decision Intelligence
- Tangible, Industry/Issue Specific AI-Applications
- Build your own apps with the Squirro AI-Studio



## Reference Customers



## Investors



# Customer Success

A self-learning system keeping you in the know & recommending what's next.

3 powerful application suites built on top of a powerful & versatile platform:

**Sales Insights:** Top line growth

**Service Insights:** Cost savings

**Risk Insights:** Reduce exposure

**Enterprise Search:** A more effective enterprise

All Solutions are SquirroGPT enabled.

Sales

CANDRIAM  
A NEW YORK LIFE INVESTMENTS COMPANY

armacell  
MAKING A DIFFERENCE AROUND THE WORLD

Sales: Source opportunities

Service

Standard Chartered

OCBC Bank

Service: Automated case resolution

Risk

EUROPEAN CENTRAL BANK  
EUROSYSTEM

DEUTSCHE BUNDESBANK  
EUROSYSTEM

Risk: Threat detection for supervisory & audit

ES

MUBADALA

Henkel

Cognitive Search: A single data & knowledge fabric



# Thank You

Dorian Selz

---

[dorian@squirro.com](mailto:dorian@squirro.com)

