

# Notes of OReilly course - Generative AI Foundations: Understand Generative AI Applications, Benefits, and Risks

Link: <https://learning.oreilly.com/course/generative-ai-foundations/0642572079628/>

What is Gen AI

---



## Summary:

- Generative AI is subset of Machine learning
- Large Language Models are trained on massive amounts of data
- There are models which can do more than just language.
- Embeddings are the way you make a LLM understand language.
- Transformer architecture
- Special training techniques like instruction tuning and RLHF
- Diffusion models are powering text to image generation

Prompt Engineering

---



## Prompt Engineering

Prompt engineering refers to the process of crafting effective and specific prompts or inputs to interact with AI models in order to get desired outputs.



## The C.R.E.A.T.E framework

Practice	What is it	Importance	What it does
<b>Character</b>	Crafting prompts with a distinct tone, persona, or specific attributes that align with your brand's identity.	Establishes a consistent brand voice, making AI-generated content feel cohesive and on-brand.	Infuses AI-generated content with a unique personality, reinforcing brand recognition and loyalty.
<b>Request</b>	Constructing clear and specific prompts that explicitly convey the desired outcome or information from the AI model.	Reduces ambiguity, ensuring the AI generates content that meets your intended purpose and requirements.	Guides the AI towards producing accurate and relevant outputs, saving time and improving efficiency.
<b>Examples</b>	Providing relevant examples within the prompts to give the AI context and a reference for generating appropriate content.	Helps the AI understand the desired output style or format, resulting in content that matches expectations.	Enhances the AI's ability to generate content that aligns with the provided examples and context.
<b>Additions</b>	Including specific instructions or keywords that guide the AI to focus on certain aspects or elements in the output.	Allows you to control the AI's creative direction and highlight specific points in the generated content.	Directs the AI's attention to particular details, ensuring the generated content is comprehensive.
<b>Type of Output</b>	Specifying the type of content or format you're seeking from the AI, whether it's a summary, list, explanation, etc.	Tailors the AI's output to match your intended presentation or communication style for better relevance.	Shapes the AI's response to match the desired content structure and format, improving communication.
<b>Extras</b>	Including additional context, relevant background information, or constraints that help the AI better understand the task.	Provides the AI with a complete picture, reducing potential confusion and leading to accurate outputs	Equips the AI with necessary information to generate content that aligns with your specific needs.



## A good prompt

You are a highly experienced writer who writes concise and readable text without stop words, filler words or jargon. I want you to summarise the following text, highlighting the most important concepts. Deliver this as a short paragraph of 100 words. Then list the most important points as a bullet-point list. Finally, follow it with a one sentence summary. The text I want you to summarise is "[TEXT]"



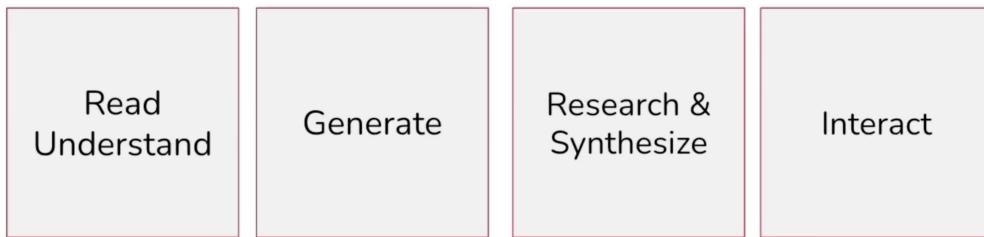
## Additional learning

<https://www.promptingguide.ai/>

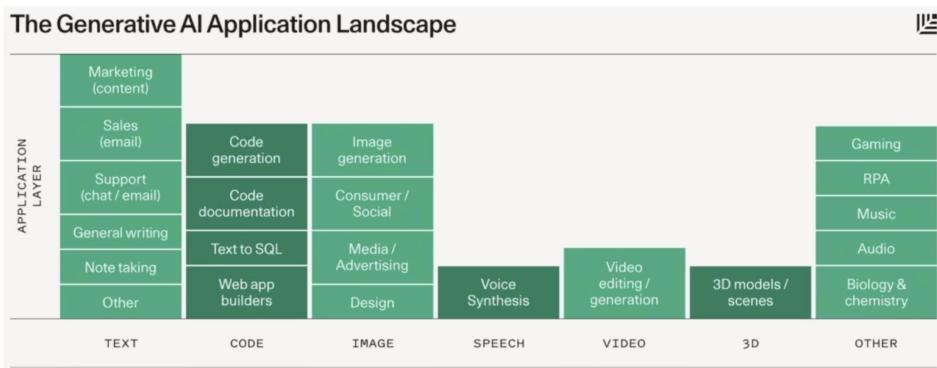
## Applications of Gen AI



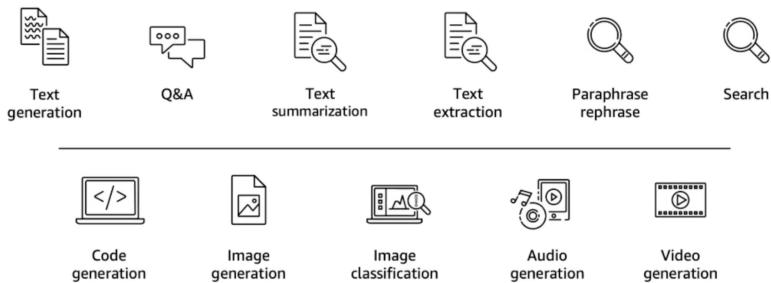
## What Generative AI does today?



## Use Cases and Applications



## Use Cases and Applications

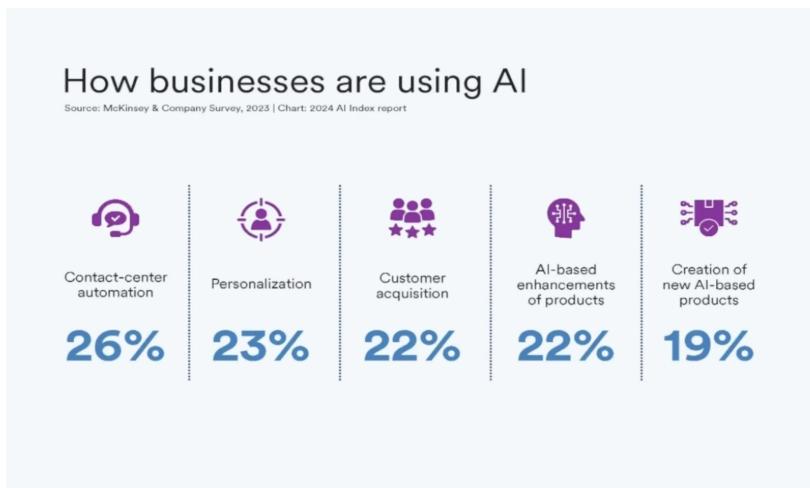




## Applications and usage



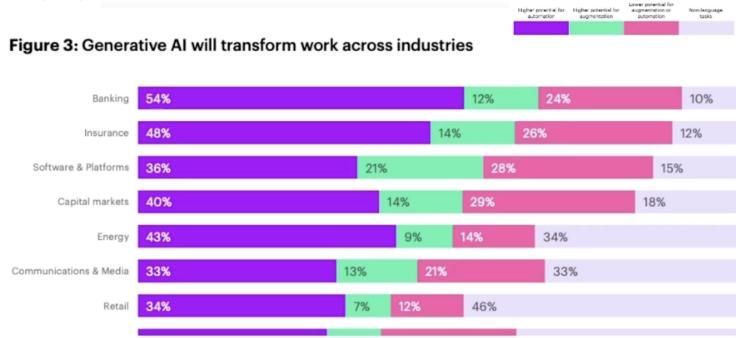
## Use Cases of GenAI



## Impact on Business

40% of working hours across industries can be impacted by Large Language Models (LLMs)

Figure 3: Generative AI will transform work across industries





## Top enterprise use cases for Gen AI

<b>Written and augmented content</b>	Producing a "draft" output of text in a desired style and length
<b>Question answering and discovery</b>	Enabling Knowledge base using private data / In-context learning
<b>Summarization</b>	Text manipulation, to soften language or professionalize text
<b>Tone</b>	Shortened versions of conversations, articles, emails and webpages
<b>Simplification</b>	Breaking down titles, creating outlines and extracting key content
<b>Classification of content</b>	Sorting by sentiment, topic, etc.
<b>Chatbot performance improvement:</b>	Sentiment classification and generation of journey flows from general descriptions
<b>Software coding:</b>	Code generation, translation, explanation and verification



## LLM use cases



## Economic Impact of Gen AI

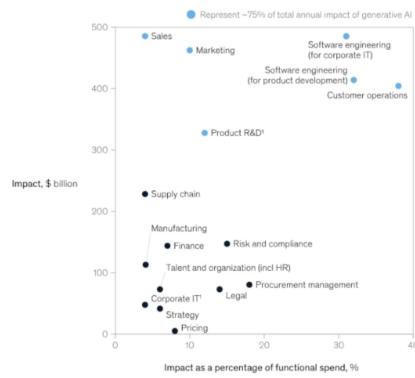
\$2.6 - \$4.4 trillion global economy

Four key areas:

1. Customer operations,
2. Marketing and sales,
3. Software engineering,
4. Research and Development (R&D).

In banking we could potentially see an extra \$200 billion to \$340 billion in value each year if we fully implement generative AI solutions.

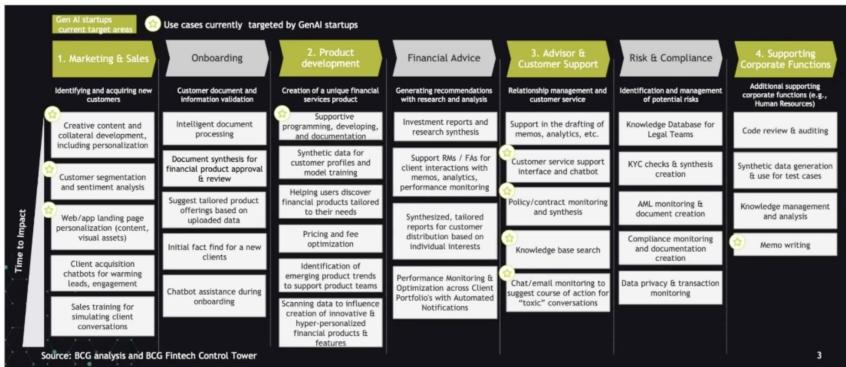
Generative AI is changing the way we work by automating tasks that currently take up 60-70% of employees' time.





# Financial services

Generative AI use cases



Financial services : Stories

01	JP Morgan	<ul style="list-style-type: none"> <li>IndexGPT to select Securities</li> </ul>
02	Goldman Sachs	<ul style="list-style-type: none"> <li>POCs, classification docs, Dev team productivity</li> </ul>
03	DBS	<ul style="list-style-type: none"> <li>Personalise product recommendations</li> </ul>
04	Alibaba	<ul style="list-style-type: none"> <li>Tongyi Qianwen, LLM</li> </ul>
05	Morgan Stanley	<ul style="list-style-type: none"> <li>Own chatbot for 16000 advisors</li> </ul>

## Limitations

## Limitations and Risks

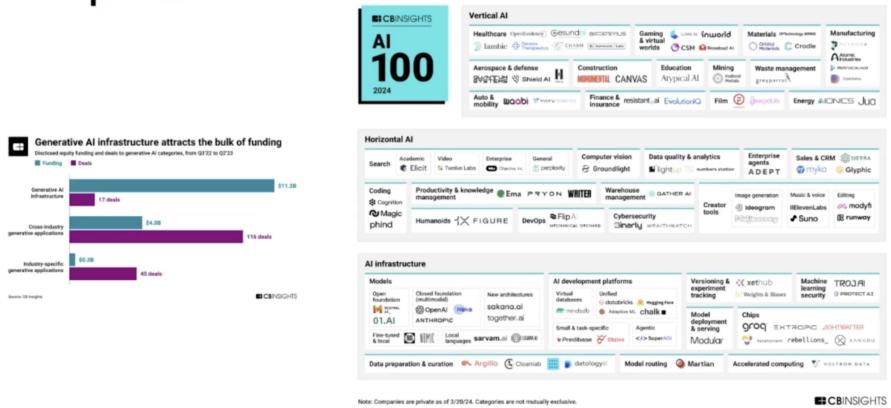
- Lack of Accuracy
- Bias and Fairness
- Misinformation / hallucinations
- Lack of Interpretability
- Ethical concerns
- Data Privacy and Security

## Responsible AI



## Companies and Tools

## Companies

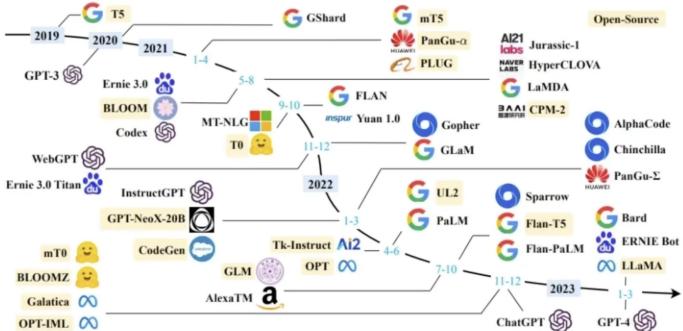


## Tools

Tool	Company	Description
Chatgpt, GPT-4	OpenAI	A large language model that can generate text, translate languages, write different kinds of creative content, and answer your questions in an informative way.
Midjourney	Midjourney	A text-to-image AI tool that can create realistic images from text descriptions.
Stable Diffusion	Stable Diffusion	A generative model that can create high-quality images from text descriptions.
GitHub Copilot	Microsoft	A tool that helps programmers write code. It is powered by OpenAI's GPT-3 language model and can suggest code completions, documentation, and other helpful information.
DALL-E 3	OpenAI	An image generation tool that can create realistic images from text descriptions.
Adobe Firefly	Adobe	A tool that can generate realistic images, videos, and 3D models from text descriptions.
D-ID	D-ID	A tool that can create realistic video avatars connected with audio and language models.
Runway.ml	Runway	A platform that allows users to create and deploy generative AI model generating videos, enhancing images from just text.

## Open vs Closed Source Models

### Open Source Models





## Open Source Models

[https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard)



## Choosing Closed vs Open Source

Closed Source (Private)	Open Source (Self hosted)
Effective for Low usage	Effective for high usage
Min Infra setup	High deployment and setup costs
High quality responses	Average quality responses
Sensitive for Privacy	Insensitive to Privacy
External dependency as a black box	Full Control over code and the model.
High Speed of implementation	Requires time and resources to implement



## Models

Open Source Models	Closed Source Models
Llama 3.1	GPT4o, GPT-4, GPT 3.5, GPT-3
Mistral	Claude
BLOOM	Cohere
BERT	Gemini
Stable Diffusion	Dall-E

Home of open source models : <https://huggingface.co/>



## Self hosting vs Vendor LLM

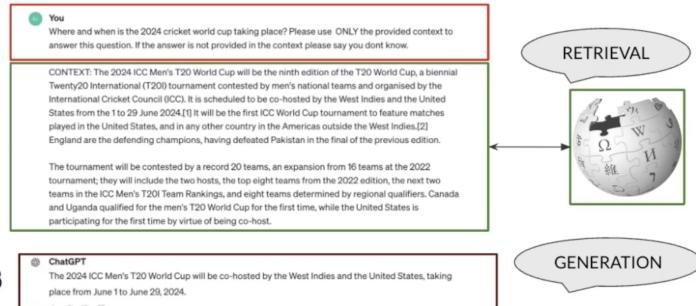
You need	You Choose
More Privacy	Self hosted LLM
Quick Development	OpenAI
Low costs	Open AI
Minimal Infra	Open AI
Full control	Self hosted LLM
No external dependencies	Self hosted LLM
Top quality	Open AI

---

RAG (Retrieval Augmented Generation)

## How it works?

ChatGPT 3.5 ~

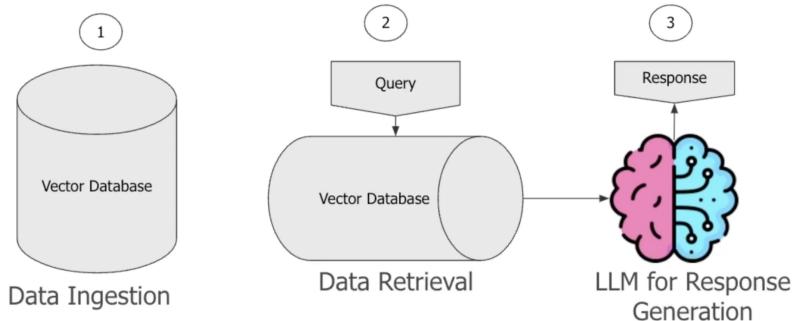


## Retrieval Augmented Generation



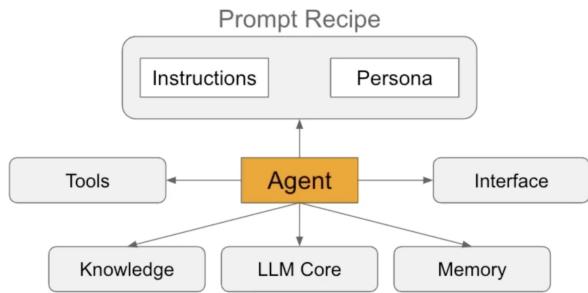
Retrieval Augmentation Generation (RAG) is an architecture that augments the capabilities of an LLM by adding an information retrieval system that provides the data.

## RAG - Overview



## Agents

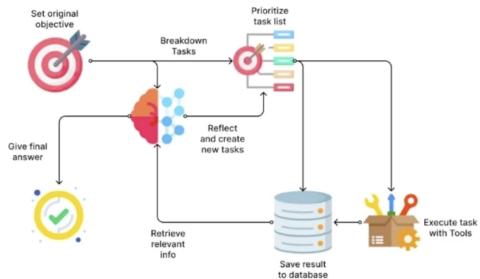
## Agents - Core concept



## Agents - BabyAGI



Apply iterative work to accomplish high level goals



## MetaGPT: Collaboration with Agents

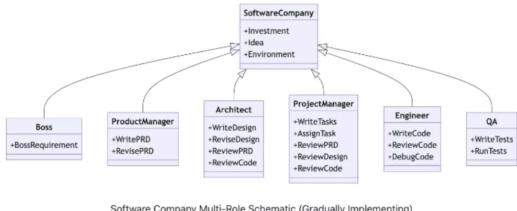


Assign different roles to GPTs to form a collaborative software entity for complex tasks.

日本語 English ドキュメント 日本語 MetaGPT ライセンス MIT ROADMAP 会社概要 WeChat 質問 GitHub Codespace Open

1. MetaGPT takes a one line requirement as input and outputs user stories / competitive analysis / requirements / data structures / APIs / documents, etc.
2. Internally, MetaGPT includes product managers / architects / project managers / engineers. It provides the entire process of a software company along with carefully orchestrated SOPs.

I. Code = SOP(Team) is the core philosophy. We materialize SOP and apply it to teams composed of LLMs.



Software Company Multi-Role Schematic (Gradually Implementing)

## MindOS: Agent marketplace in action



The screenshot shows a user interface for a job application. On the left, there's a sidebar with 'Marketplace' and a list of agents: 'My Geniuses' (Lupe, Jeff), 'Ksoon', 'Chris', 'Hans', and 'Dante Friend'. Below that is a 'Create New Service' button. The main area has tabs for 'Subscription' and 'Ksoon'. A large video player in the center displays a video of a woman in a black blazer. The video content is a self-introduction for a Product Hunt profile. The video player includes controls like play/pause, volume, and a progress bar showing 1:14 / 1:34. At the bottom right is a 'YouTube' icon.

## Collaboration with Agents



### Generative Agents: Interactive Simulacra of Human Behavior

03442v2 [cs.HC] 6 Aug 2023

The image shows a top-down map of a campus or town layout. Overlaid on the map are several small windows representing agent interactions:

- Taking a walk in the park:** Shows a person walking through a park area.
- Joining for coffee at a cafe:** Shows two people at a cafe counter. One says, "May I join you for coffee? I have some news to share. How are you?"
- Arriving at school:** Shows a person entering a school building.
- Sharing news with colleagues:** Shows two people talking. One says, "Linda! May I tell you about the upcoming agency acquisition?"
- Finishing a morning routine:** Shows a person in a bathroom.

Below the map, there is a legend with icons for a person, a house, a car, and a bus, each associated with a different color: blue, green, red, and yellow respectively.

Some Agents

The screenshot shows the AgentGPT Beta interface. On the left, there's a sidebar with a 'Rework' section containing multiple 'TravelGPT' entries and a 'Pages' section with 'Templates', 'Help', 'Settings', and 'Manage account'. Below these are social sharing buttons for Twitter, LinkedIn, and a 'Subscribe' button. A user profile for 'Altaf Rehmani' is also present. The main area features a header 'AgentGPT Beta' with a 'Beta' badge. Below it is a banner: 'Interested in AI Agents to scrape web data? Find out more here >'. A message encourages users to 'Create an agent by adding a name / goal, and hitting deploy! Try our examples below!'. Three AI agent cards are displayed: 'ResearchGPT' (Create a comprehensive report of the Nike company), 'TravelGPT' (Plan a detailed trip to Hawaii.), and 'StudyGPT' (Create a study plan for a History 101 exam about world events in the 1980s). At the bottom, there's a form to 'Name' the agent as 'AgentGPT' and set the 'Goal' to 'Make the world a better place', with play, pause, and stop buttons.

The screenshot shows the MindOS Studio Marketplace page. It has a header 'MindOS STUDIO' and a 'Marketplace' tab. The 'Featured AI Minds' section displays three cards: 'Industry Analyst' (Assist you in conducting desktop research, analyze an in...), 'Trip Advisor' (Help you plan your trip and recommend great places), and 'News Genie' (News roundup and deep dive). Below this is a 'Trending AI Minds' section with a grid of AI mind icons. A video player at the bottom shows a video titled 'analogical.com' with a timestamp of 03:06 / 04:15.

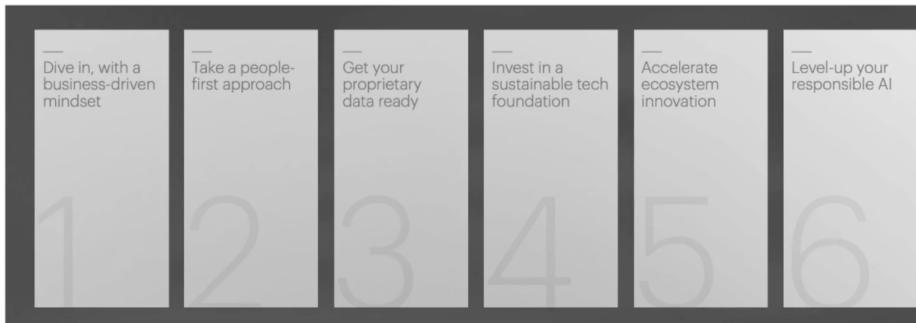
The screenshot shows a YouTube video player for a video titled 'Pro Runway Secrets + Elon Takes on MidJourney' by 'Curious Refuge'. The video has 160K subscribers and 315 likes. The player includes standard controls like play/pause, volume, and progress bar. Overlaid on the video are several AI annotation tools: 'Transcript', 'Summary', 'Notes', 'AI Chat', and 'Basic'. The 'Summary' tool shows a mind map with a central node 'AI Film News of the Week' connected to other nodes. The 'AI Chat' tool has a message from 'Elon Musk' about launching an innovative image generator on X. The overall interface is designed to facilitate analysis and interaction with the video content.

## How to get started / Key Take Aways



## How to get Started

Adoption Essentials



## Key Takeaways

- LLMs are next word predictors based on networks which have “attention” and are autoregressive
- Move towards Factual Information, built into everyday experiences
- Wide range of Applications in every industry
- MultiModal Models - Images, Audio, Video, 3D
- Rise of Open Source Models & RAG
- Responsible AI to mitigate Risks and Limitations
- Agents: The Future of Autonomous Tasks



## Future

Are you ready?

Before long, GenAI will greatly impact product development, customer experience, employee productivity and innovation. We predict that:

	By 2025, 70% of enterprises will identify the sustainable and ethical use of AI among their top concerns.		By 2025, 30% of outbound marketing messages from large organizations will be synthetically generated. That's up from less than 2% in 2022.		By 2030, AI could reduce global CO2 emissions by 5 to 15% and consume up to 3.5% of the world's electricity.
	By 2025, 35% of large organizations will have a chief AI officer who reports to the CEO or COO.		Through 2026, despite all the advancements in AI, the impact on global jobs will be neutral — there will not be a net decrease or increase.		By 2030, decisions made by AI agents without human oversight will cause \$100 billion in losses from asset damage.
	By 2025, the use of synthetic data will reduce the volume of real data needed for machine learning by 70%.		By 2033, AI solutions will result in more than half a billion net-new human jobs.		