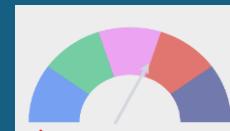


GEN AI LEARNING

Github Repository link --<https://github.com/Abhishekrai129>



\$13.7 billion in 2023
to **\$165 billion** by
2032.



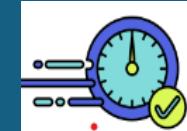
CAGR: **40.2%** from 2021 to
2028.



75% of customers prefer AI
solutions that integrate
easily with existing system



60% are concerned
about data privacy and
security.



65% of customers expect
solutions to provide
instant responses



80% of users
prefer tailored
recommendations

Value Creation and Monetization in Gen AI

VALUE PROPOSITIONS

Data-Driven Insights: 70% improve decision-making (McKinsey)

Increased Efficiency: 40% time saved (Harvard Business Review)

Cost Reduction: Up to 30% lower costs (Deloitte)

Personalized Experiences: 80% customer satisfaction (Epsilon)

Scalability: 50% faster growth (Gartner)

CUSTOMER SEGMENTS

Enterprises: 85% using AI (PwC)

SMBs: 40% adopting AI tools (Forrester)

Developers: 60% prefer AI integrations (Stack Overflow)

Startups: 70% leverage AI for growth (TechCrunch)

Public Sector: 50% using AI solutions (Accenture)

REVENUE STREAMS

Subscription Fees: \$300 billion market (Statista)

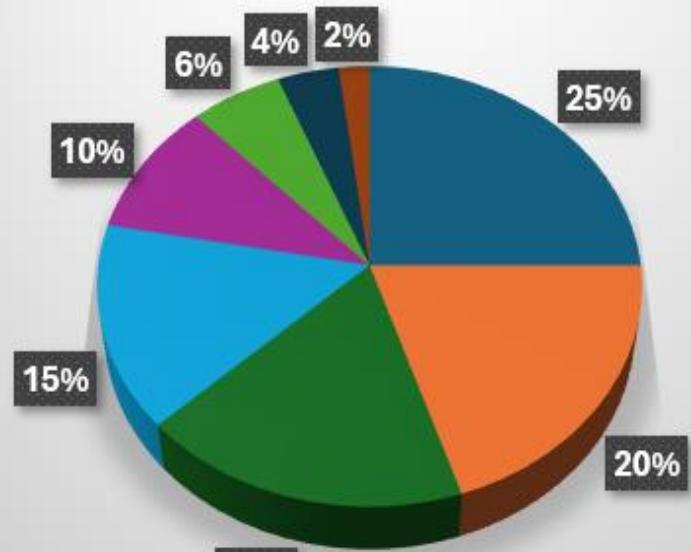
Licensing: 20% revenue from licensing (Bain & Company)

Consulting Services: \$200 billion industry (IBISWorld)

Training Programs: 60% of clients prefer training (Learning Guild)

Data Monetization: \$1 trillion potential (McKinsey)

INDUSTRY TRENDS IN GEN AI



Healthcare

Retail

Telecommunications

Media & Entertainment

Finance

Manufacturing

Transportation & Logistics

Other Industries

The Building Blocks of Generative AI

EXPANDING CAPITAL

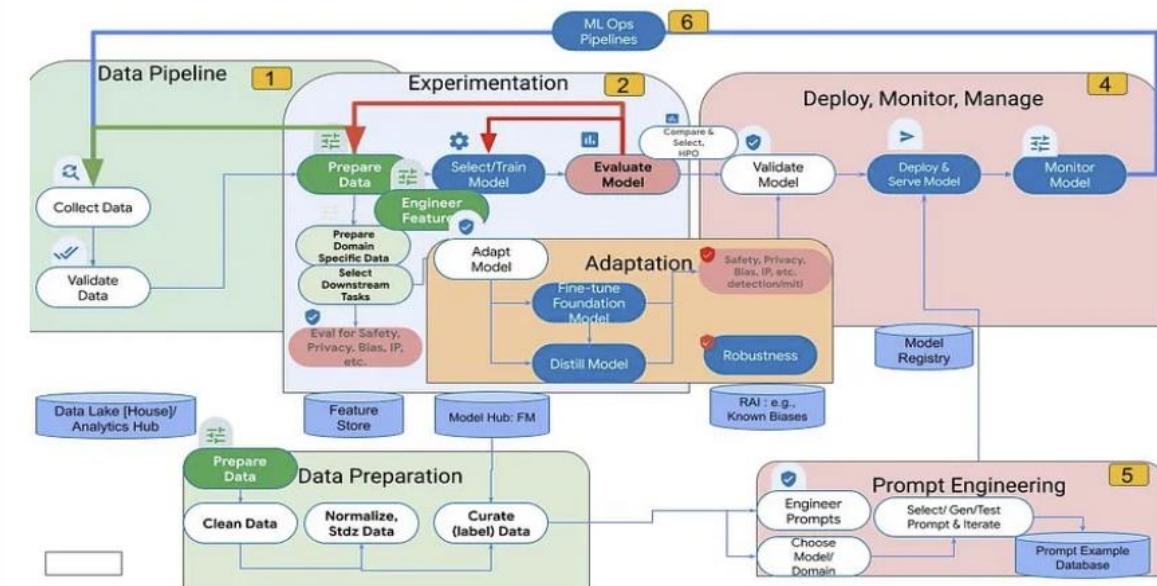
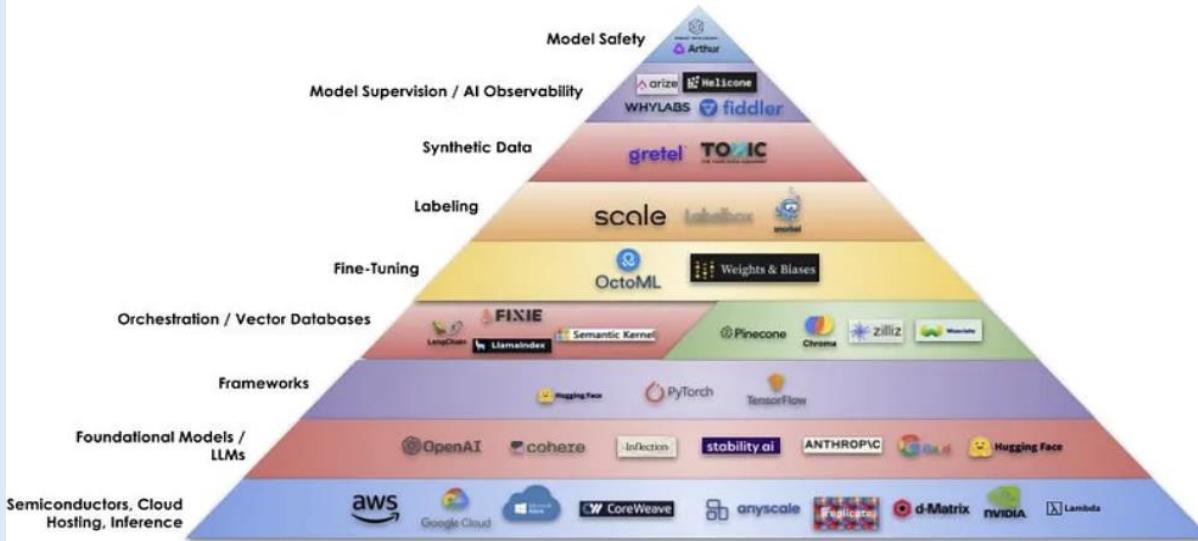
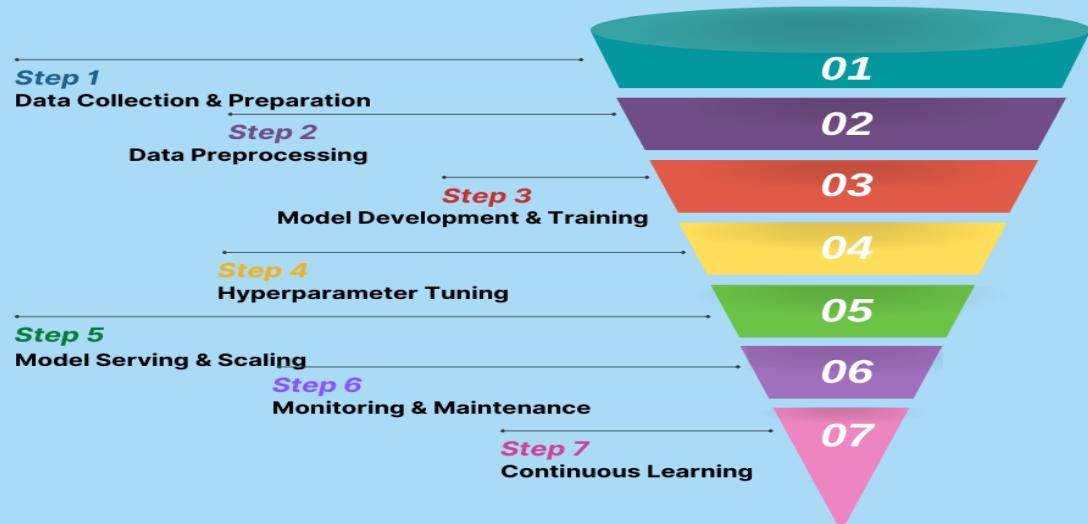


Figure 1: The Generative AI Lifecycle

GEN AI TECH STACK



Data Collection-Gather and store datasets.

Tools: Scrapy, Twitter API, AWS S3, Snowflake

Data Preprocessing-Clean and prepare data.

Tools: Pandas, PySpark, SpaCy, NLTK

Model Development-Train machine learning models.

Tools: TensorFlow, PyTorch, Hugging Face Transformers,

Hyperparameter Tuning-Optimize performance.

Tools: Optuna, Ray Tune, Hyperopt

Model Serving-Deploy models.

Tools: AWS SageMaker, Docker, TensorFlow Serving, Kubernetes

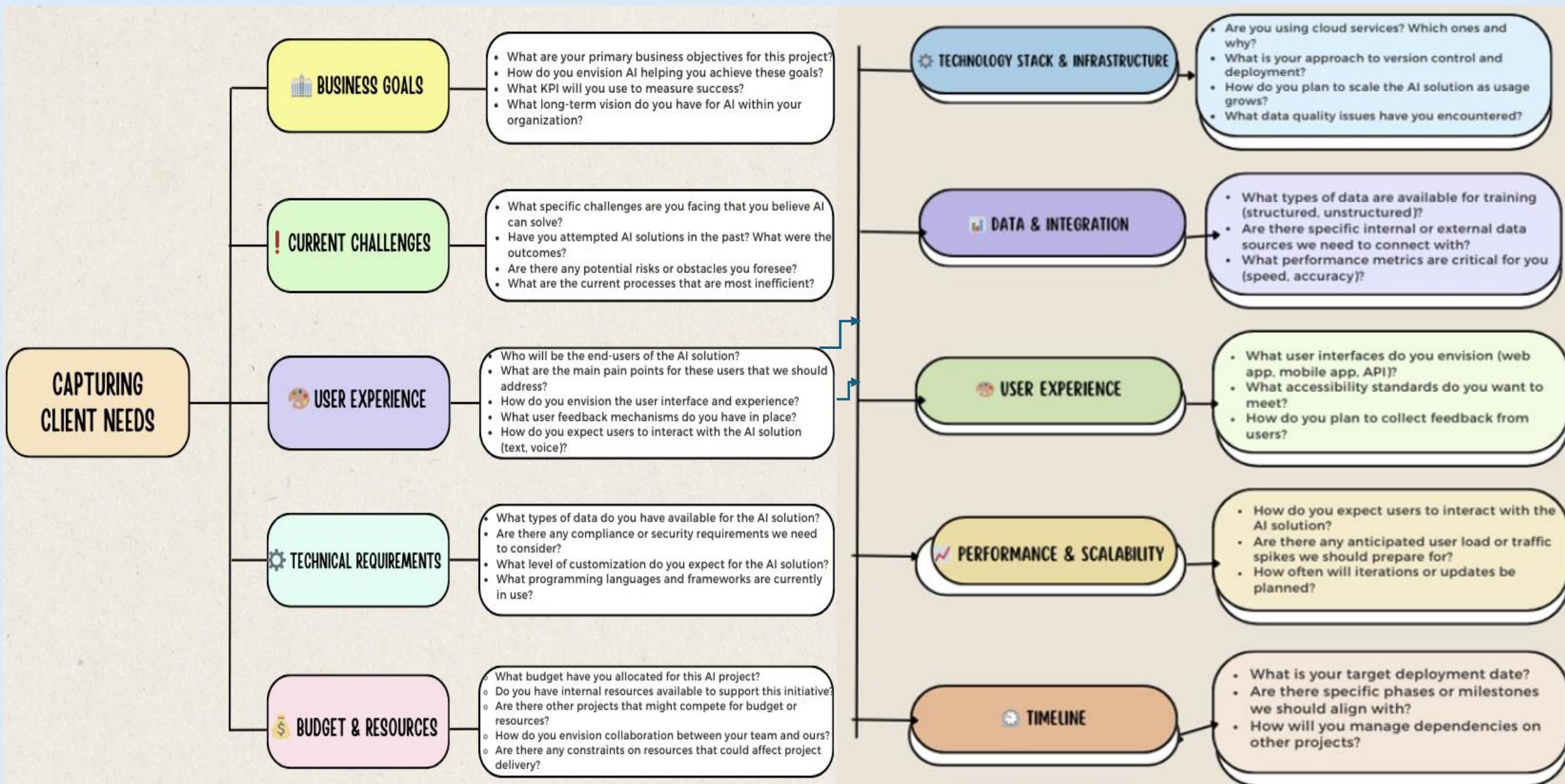
Monitoring=Track model performance.

Tools: Prometheus, Grafana, ELK Stack, Apache Kafka

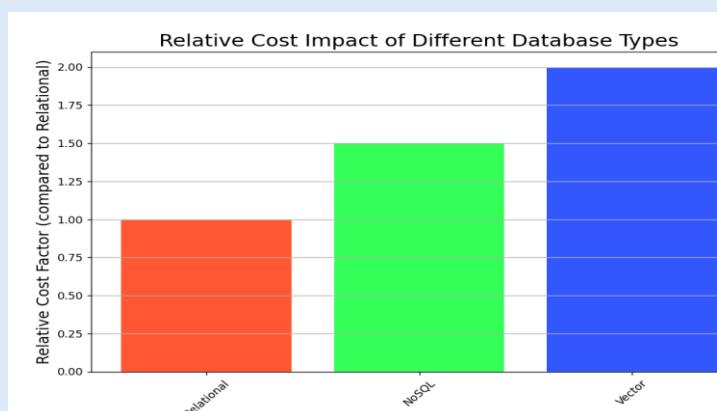
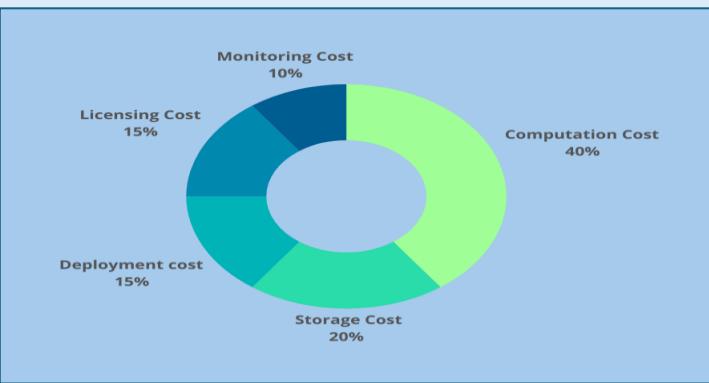
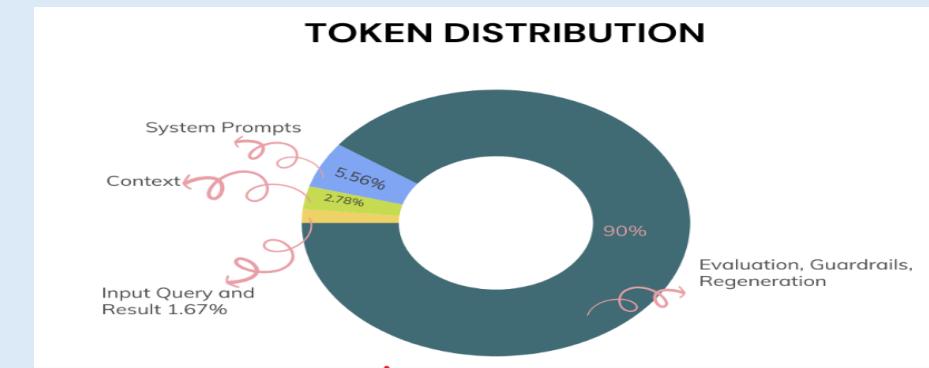
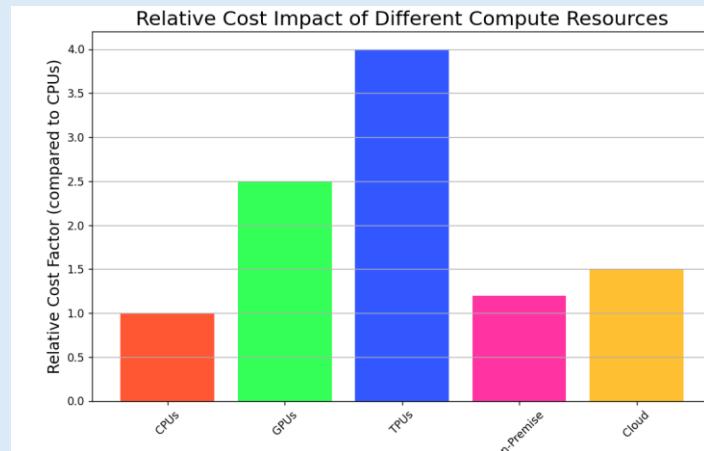
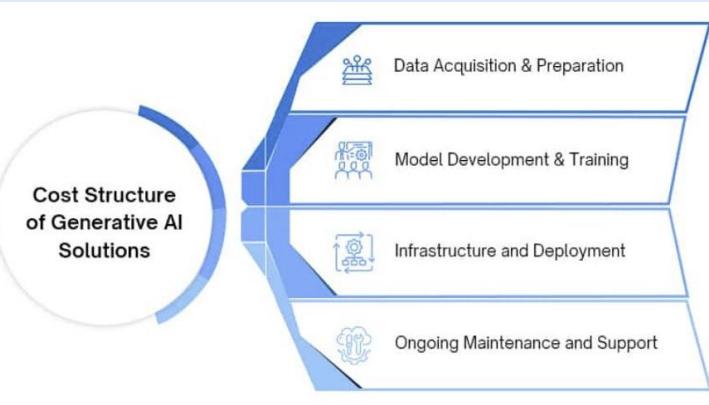
Continuous Learning=Improve models continuously with new data.

Tools: MLflow, DVC, A/B Testing (Optimizely)

Understanding Client Needs for Generative AI Projects



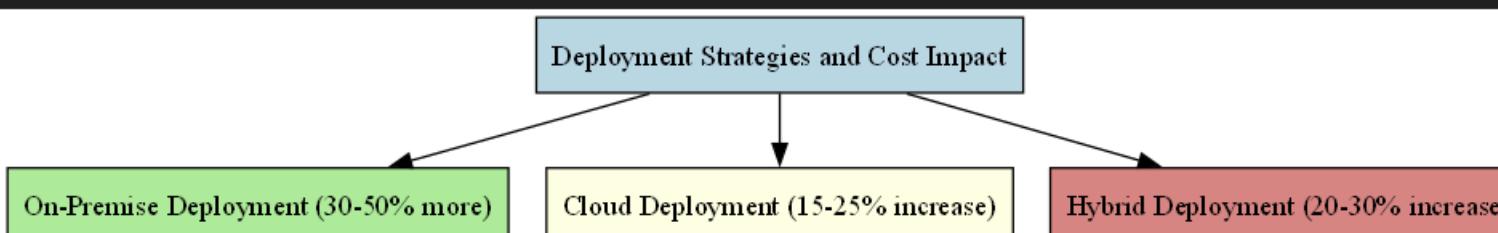
GEN AI Project COST



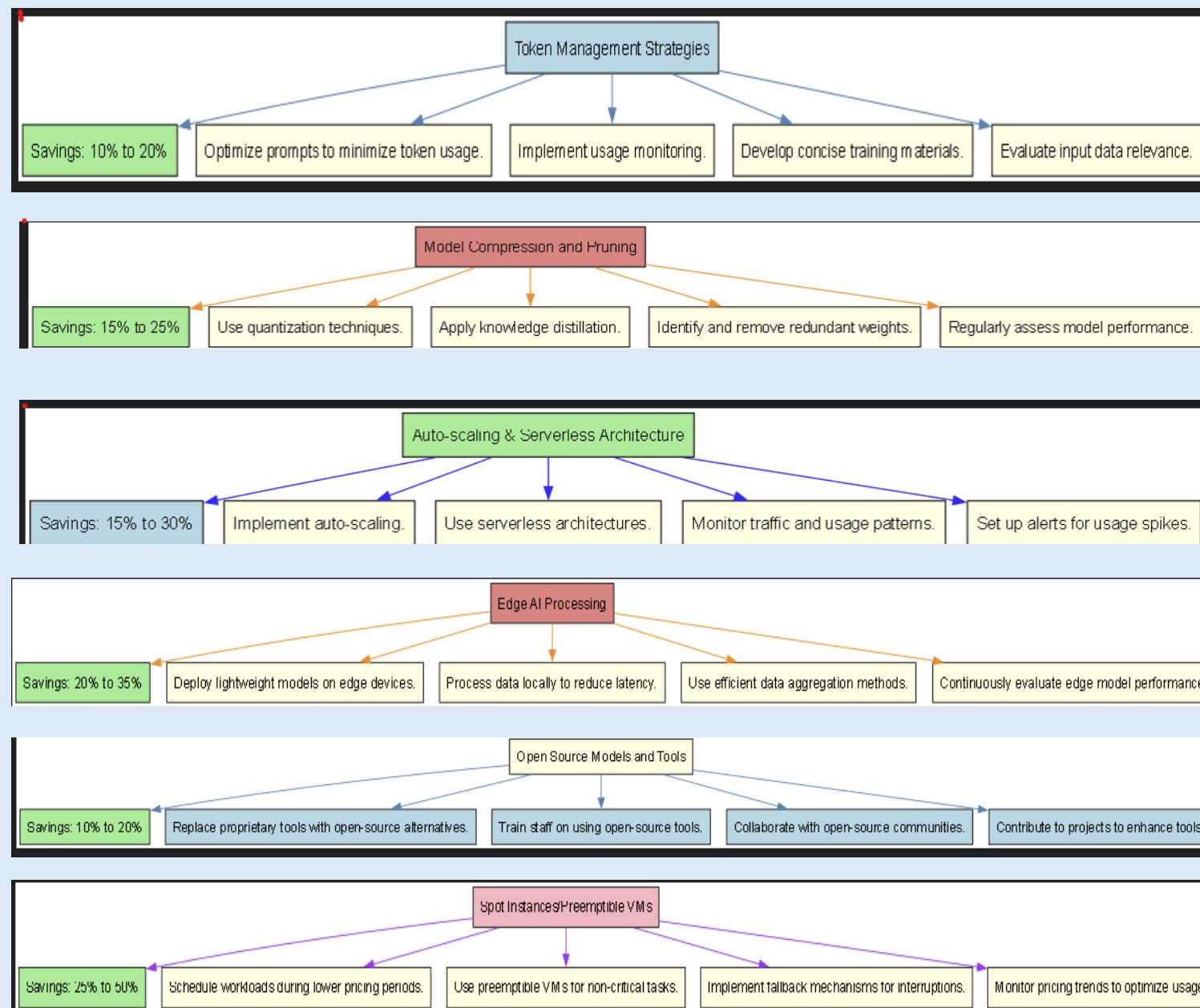
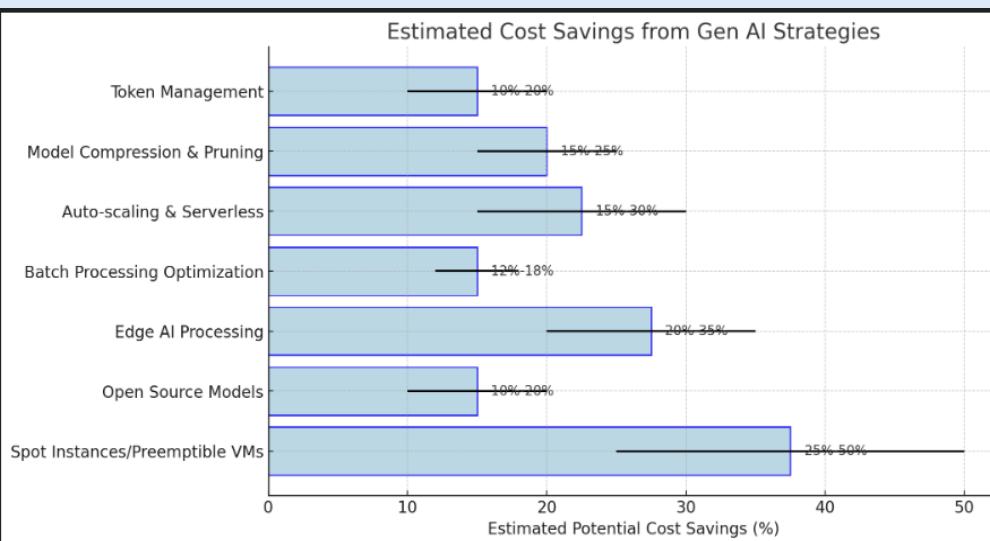
Component	Description	Example Tokens
Input/Output Tokens	Basic user queries and AI responses.	300
Context Tokens	Reference previous interactions for better responses.	500
Metadata/System Prompts	User info and instructions for task accuracy.	600
Evaluation Tokens	Guardrails for quality and safety.	100
Error Handling/Regeneration	Tokens for retries or clarification.	300+
Total Tokens per Interaction	Sum of all token costs.	1,800
Annual Cost (10 million interactions)	Cost for GPT-4 example.	\$465,000

Available GPUs (by cloud provider)

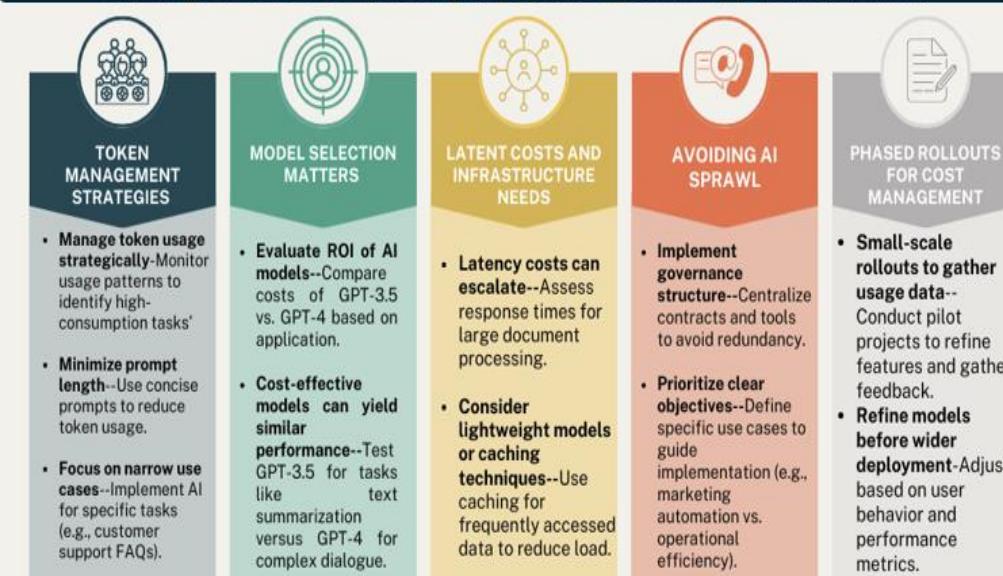
Cloud Service Provider (CSP)	Available GPUs	Pricing Indicator (A100 40gb, on-demand)
AWS	A100, V100, M60, T4, A10, Trainium, Inferentia, Gaudi, V520	\$4.10/h (p4d.24xlarge)
Azure	H100, A100, V100, P100, K80, M60, P40, T4, A10, MI25	\$3.40/h (ND96asr A100 v4)
Oracle	A100, V100, P100, A10	\$3.05/h (BM.GPU4.8)
Google	A100, V100, P100, K80, T4, P4, TPUV4	\$2.93/h (NVIDIA A100 40GB)
Coreweave	H100, A100, V100, A40, A6000, A5000, A4000, Quadro RTX 4k/5k	\$2.06/h (A100 40GB NVLINK)
Lambda Labs	A100, V100, A6000, A10, Quadro RTX 6k	\$1.10/h (tx NVIDIA A100)
FluidStack	A100, V100, A40, A6000, A5000, QUADRO RTX 5k/4k, RTX 3090, RTX 3080, RTX 2080Ti	\$1.73/h (A100 40 GB)



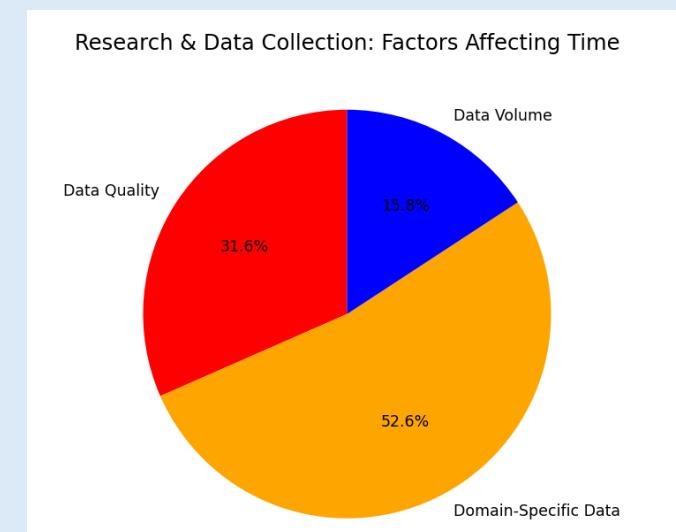
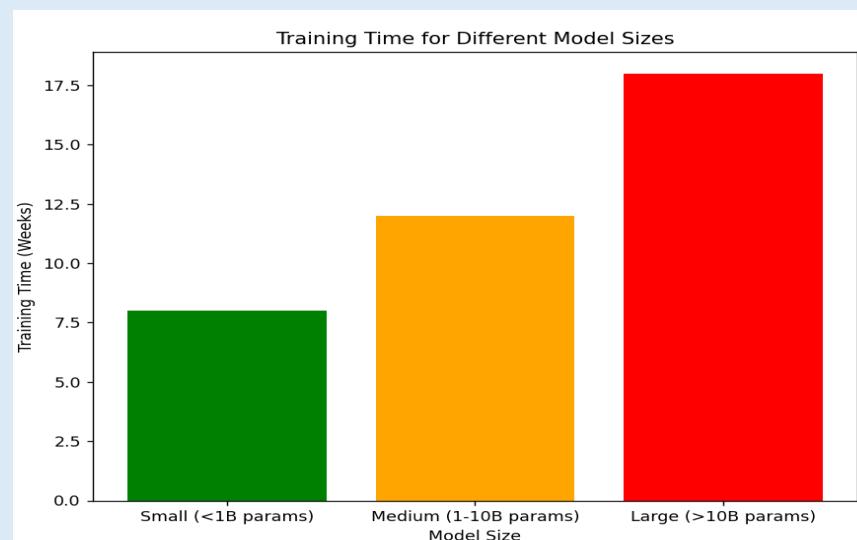
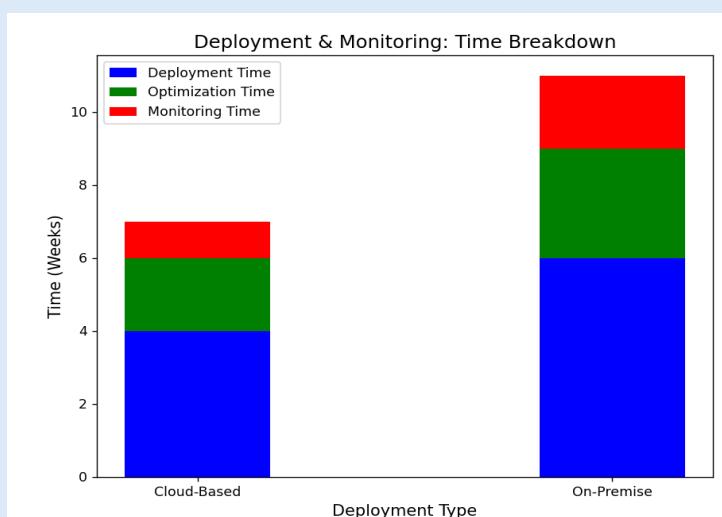
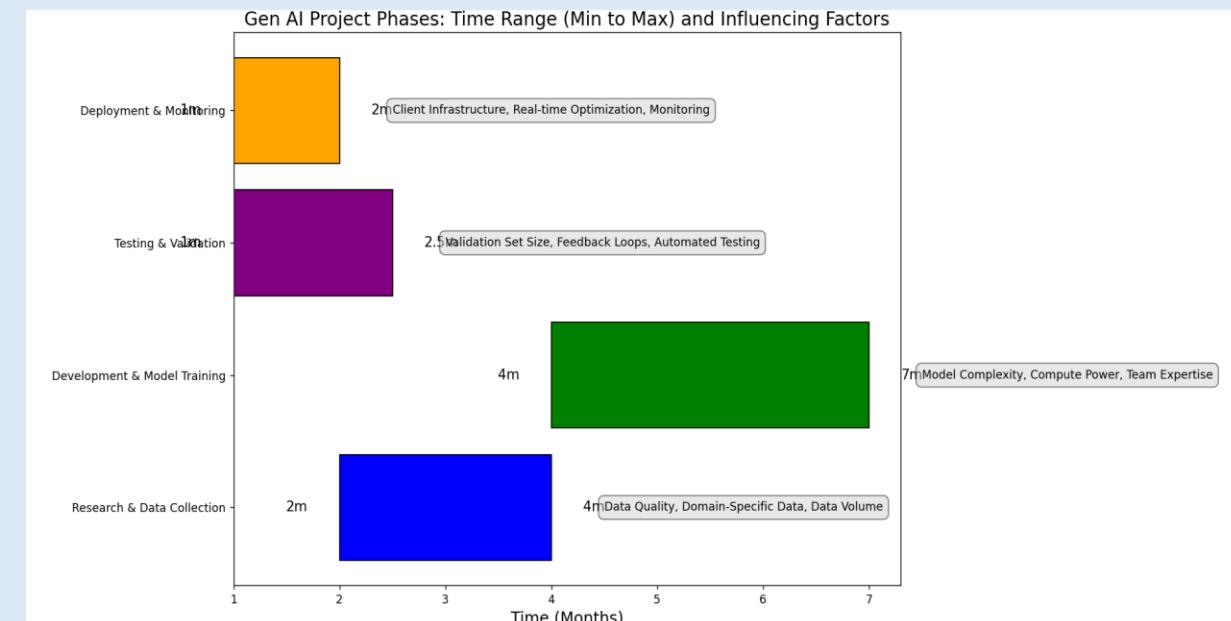
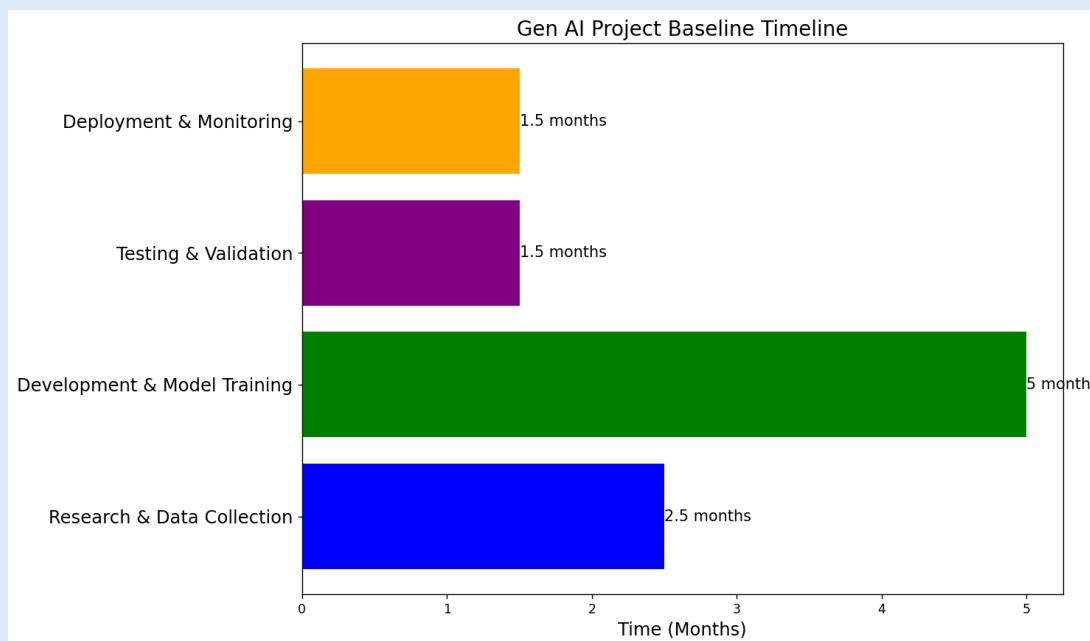
Strategies for optimizing Costs



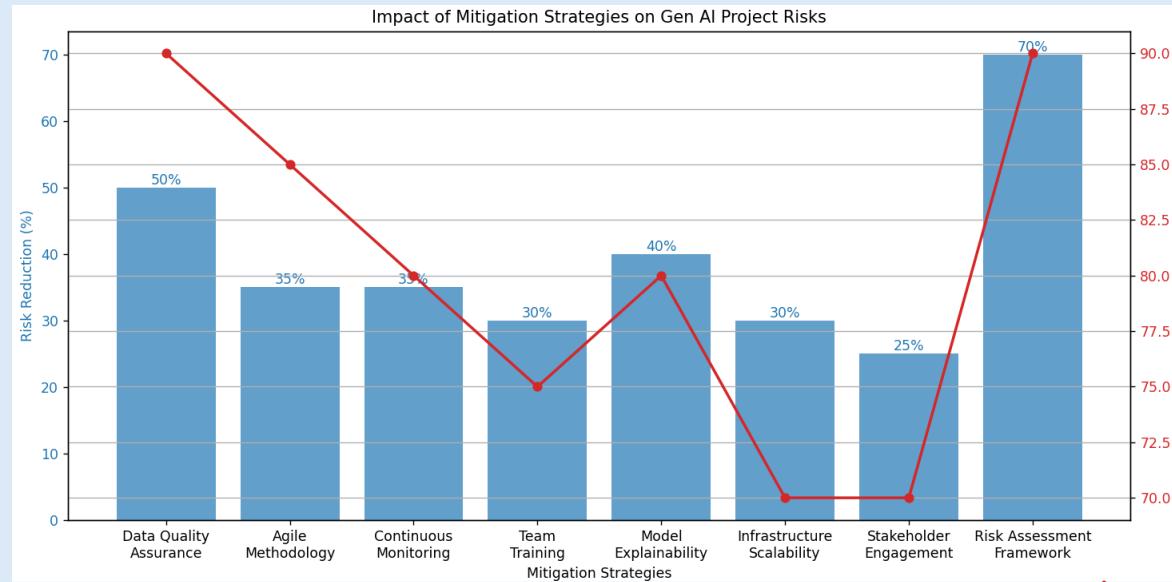
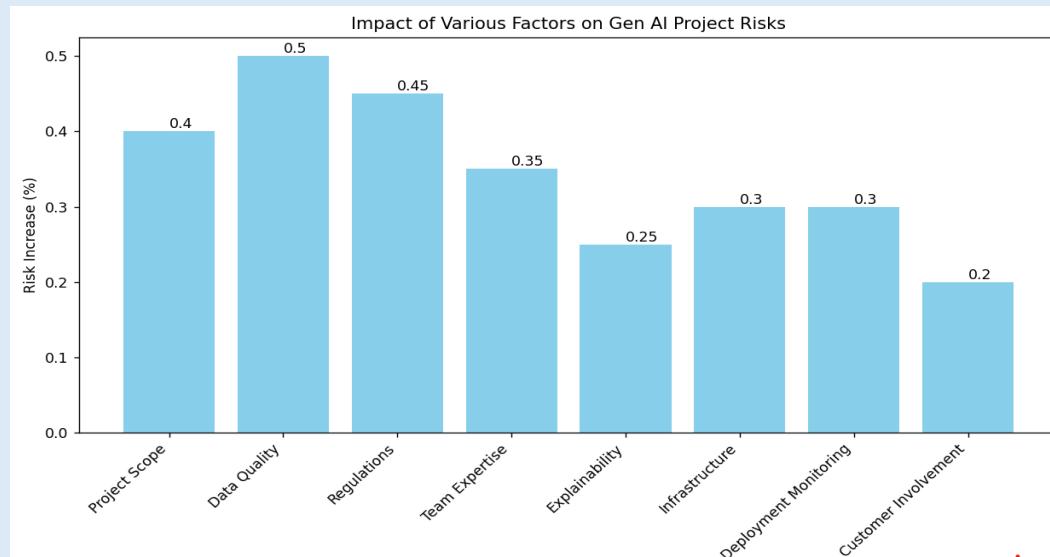
Strategies for Cost Management in GEN AI Implementation



GEN AI SCHEDULE



GEN AI RISK



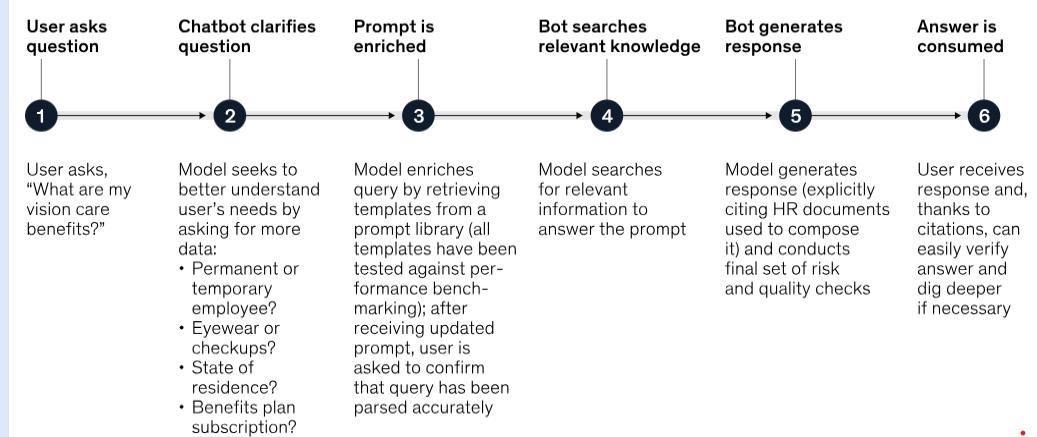
Different generative AI use cases are associated with different kinds of risk.

Generative AI use case	Impaired fairness	IP ¹ infringement	Data privacy and quality	Malicious use	Security threats	Performance and 'explainability'	Strategic
Customer journeys (eg, chatbots for customer services)	✓		✓			✓	✓
Concision (eg, generating content summaries)	✓	✓				✓	
Coding (eg, generating or debugging code)		✓		✓	✓	✓	
Creative content (eg, developing marketing content)	✓	✓		✓		✓	

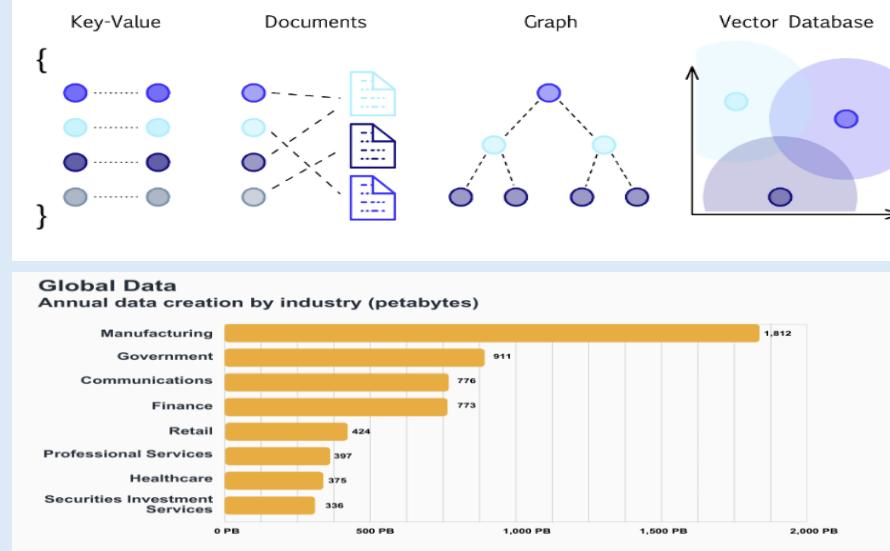
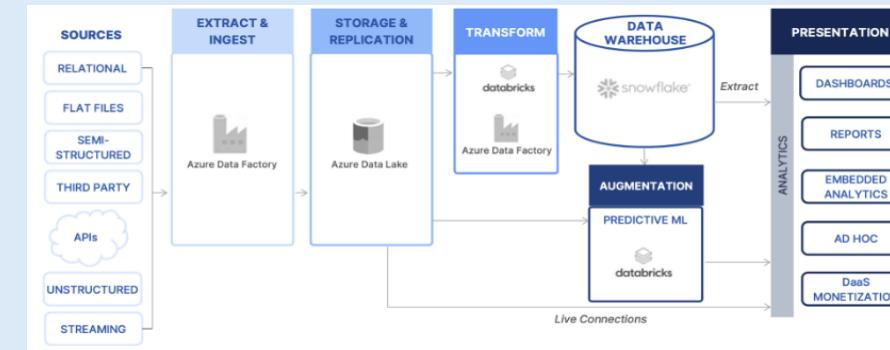
✓ Primary risk

Generative AI risk can be mitigated at multiple points across a user interaction.

Sample HR chatbot interaction with built-in checkpoints to catch potential misfires



DATABASE



Critical capability	Why this matters	SingleStore Cloud OLTP+OLAP vector capable	Pinecone vector database	Weaviate vector database	Milvus vector database	MongoDB Atlas NoSQL OLTP vector capable	Cosmos DB NoSQL vector capable	ElasticSearch full-text search
SPEED	Immediate insights; Responsive application experience for end users							
In-memory performance	Responsive applications and instantaneous analytics	●	○	○	●	○	○	○
Streaming data ingestion	Immediate queryability of streaming data from multiple sources (files, Kafka, Spark, HDFS or object stores such as S3)	●	○	○	●	●	●	○
Columnstore	Low-latency (~10s of milliseconds) on complex queries	●	○	○	●	○	○	○
Analytics on semistructured data	High-performance analytics on relational and semi-structured data in the same database engine	●	○	○	●	●	○	●
SCALE	Adapt to growing needs							
Horizontal scalability	Distributed (shared-nothing) architecture that decouples storage and compute to allow scaling using low-cost infrastructure	●	●	●	●	●	●	●
Read replicas/multi-region deployments	Run multiple workloads and scale compute instances across shared databases	●	●	●	●	●	●	●
Resiliency	Run critical applications and workloads; Mitigate risks on business operations and reputation	●	●	●	●	●	●	○
Deploy anywhere	Ability to deploy both as a fully managed cloud service or self-managed on-premises	●	○	●	●	●	○	●
SIMPLICITY	Minimize complexity and costs							
SQL-powered OLTP + OLAP with Zero ETL	Minimize data movement and duplication; minimize complexity and costs emanating from sprawl; power and simplicity of SQL for CRUD and rich query operations.	●	○	○	○	○	○	●
Multi-model	Ability to store and query multiple secondary data formats (full-text search, JSON, time series, geospatial, etc.)	●	○	○	○	●	●	○
Vector search engine	Efficiently handle large amounts of vector data for vector similarity search	●	●	●	●	●	●	●
Open-source software	Community-developed software that's typically free to use and distribute	○	○	●	●	●	●	●

1. Pure Vector Databases



2. Full-text search databases

This category includes databases such as Elastic/Lucene, OpenSearch and Solr.



4. Vector-capable NoSQL databases



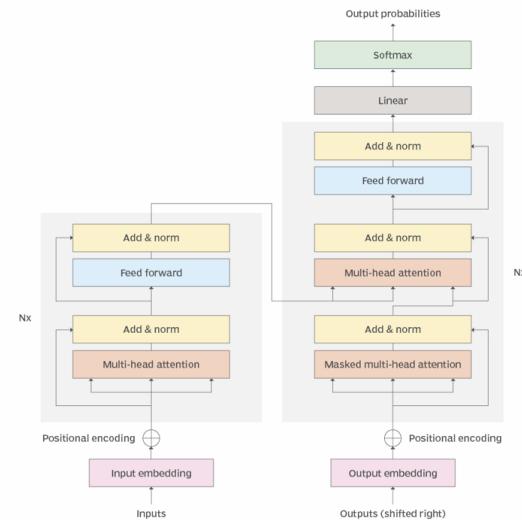
5. Vector-capable SQL databases



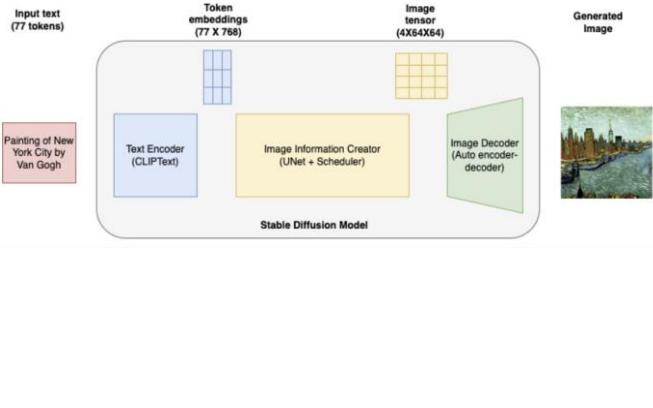
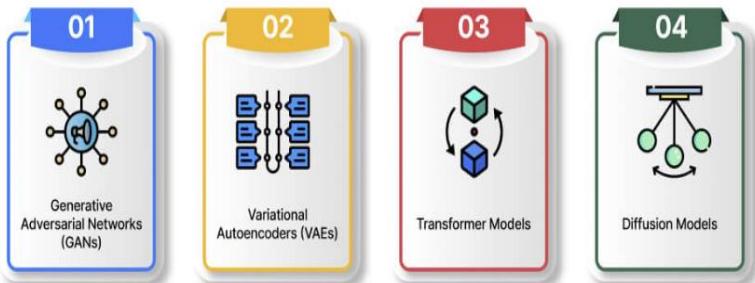
GEN AI KEY METRICS

Model Type	Success Metrics	LLM	Computer Vision	Image or Content Gen.
1. NLP & LLM (Natural Language Processing)	<ul style="list-style-type: none">- Accuracy- Precision- Recall- F1 Score- BLEU Score- ROUGE Score- Perplexity- AUC-ROC	<h3>Requests Per Minute</h3> <p>Maximum concurrent requests the service can handle per minute. Critical for large workloads.</p>	<h3>Frames Per Second (FPS)</h3> <p>Measures the number of frames the system can process per second. Essential for real-time video processing.</p>	<h3>Content Quality</h3> <p>Often assessed through human evaluations or automated metrics like BLEU or ROUGE (for text), or Inception Score and FID (for images).</p>
2. Computer Vision	<ul style="list-style-type: none">- Mean Average Precision (mAP)- Intersection over Union (IoU)- F1 Score- Accuracy- Recall- Precision	<h3>Time To First Token (TTFT)</h3> <p>Time before the service returns the first token. Important for interactive applications.</p>	<h3>Inference Time</h3> <p>The time taken to run inference on a given image. Important for applications needing fast results.</p>	<h3>Generation Speed</h3> <p>Time taken to generate an image or a text segment. Impacts usability in time-sensitive applications.</p>
3. Audio and Speech Processing	<ul style="list-style-type: none">- Word Error Rate (WER)- Signal-to-Noise Ratio (SNR)- F1 Score- Precision and Recall- Mel-Frequency Cepstral Coefficients (MFCCs)	<h3>Inter-Token Latency (ITL)</h3> <p>Average time between consecutive tokens. Affects responsiveness.</p>	<h3>Model Accuracy</h3> <p>Evaluates the correctness of the model's predictions. Often measured using metrics like precision, recall, and F1-score.</p>	<h3>Model Creativity</h3> <p>Measures the diversity and novelty of outputs. Evaluated through user studies or variety indices.</p>
4. Generative Models (GANs, VAEs, etc.)	<ul style="list-style-type: none">- Inception Score (IS)- Fréchet Inception Distance (FID)- Mean Squared Error (MSE)- Visual Turing Test	<h3>End-to-End Latency</h3> <p>Total time to complete a request. Key overall speed measure.</p>	<h3>Throughput</h3> <p>The number of images processed per second. Critical for high-throughput applications like surveillance.</p>	<h3>Resource Utilization</h3> <p>The computational resources required during generation, including CPU, GPU, and memory.</p>
5. Reinforcement Learning & Robotics	<ul style="list-style-type: none">- Cumulative Reward- Average Reward per Episode- Success Rate- Learning Curve- Task Completion Rate	<h3>Cost Per Request</h3> <p>Critical to assess cost-effectiveness for production usage.</p>	<h3>Cost Per Inference</h3> <p>Calculating the cost-effectiveness per image processed. Important for budgeting and scaling.</p>	<h3>Cost Per Generation</h3> <p>Analyzes the cost-effectiveness of generating each instance, crucial for large-scale deployment.</p>

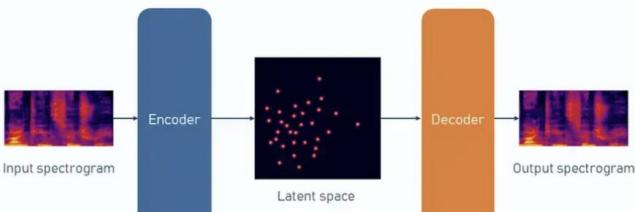
Transformer model architecture



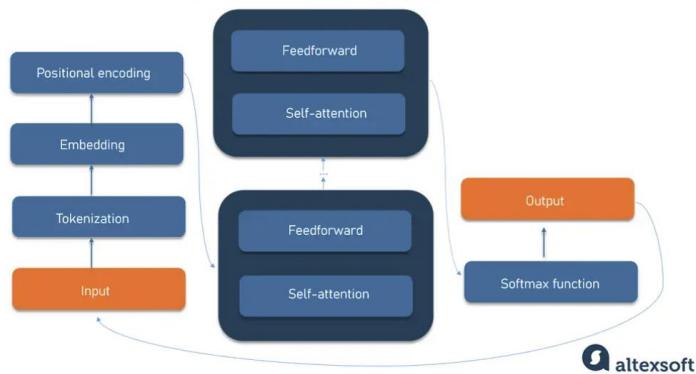
Types of Generative AI Models



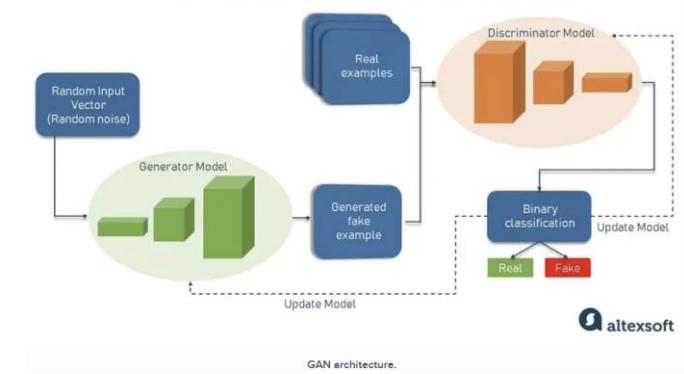
VARIATIONAL AUTOENCODERS: INCODING & DECODING LOGIC



TRANSFORMER-BASED ARCHITECTURE



GENERATIVE ADVERSARIAL NETWORK ARCHITECTURE

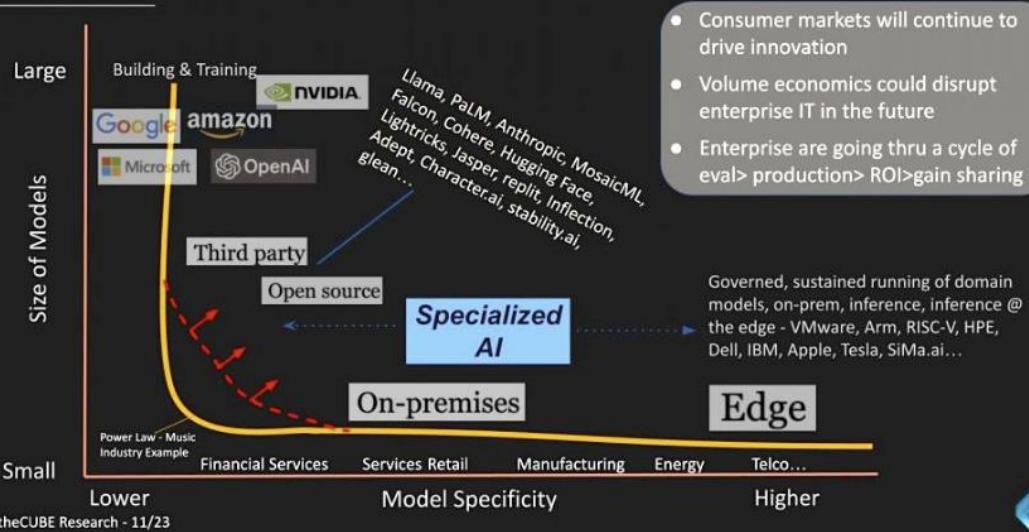


altexsoft

General AI Models and Architectures Overview				
Icon	Architecture	Description	Notable Models/Variants	Use Cases
	Transformer Architecture	Designed for sequence-to-sequence tasks; backbone of state-of-the-art models in NLP and beyond.	BERT, GPT-2, GPT-3, T5, BART	Language translation, text summarization
	Generative Adversarial Networks (GANs)	Consist of two neural networks (generator and discriminator) that compete, enabling realistic data generation.	StyleGAN, CycleGAN, Pix2Pix, cGAN	Image generation, video synthesis
	Variational Autoencoders (VAEs)	Type of autoencoder that learns data representation in latent space and enables data generation by sampling from this space.	-	Image generation, anomaly detection
	Diffusion Models	Generate data by gradually transforming a simple distribution into complex data distributions through diffusion steps.	DALL-E 2, Stable Diffusion	Image generation, denoising
	Recurrent Neural Networks (RNNs)	Used for generating sequences, especially in time-series forecasting and text generation, although somewhat overshadowed by Transformers.	LSTM, GRU	Speech recognition, text generation
	Multi-Modal Architectures	Process and generate data across multiple modalities (e.g., text, images, audio) simultaneously.	CLIP, DALL-E	Content generation, multi-modal learning
	Reinforcement Learning from Human Feedback (RLHF)	Involves training models based on human feedback, enhancing the ability to generate desired outputs.	-	Dialogue systems, content generation
Model		Key Tasks	Frameworks/Libraries	Platforms
Transformers		NLP: text classification, language modeling, machine translation	Hugging Face Transformers, PyTorch, TensorFlow	Hugging Face Platform
GPT		Natural language generation, text summarization, dialogue systems	OpenAI's GPT models, Hugging Face Transformers	Hugging Face Platform, Langchain
BERT		NLP: question answering, named entity recognition, sentiment analysis	Hugging Face Transformers, PyTorch, TensorFlow	Hugging Face Platform
GANs		Image generation, text generation, data augmentation	PyTorch, TensorFlow, Keras	-
ResNet		Computer vision: image classification, object detection	PyTorch, TensorFlow, Keras	-
YOLO		Real-time object detection	PyTorch, TensorFlow, Darknet	-
RNNs		Sequence modeling, speech recognition	PyTorch, TensorFlow, Keras	-
VAEs		Generative modeling, dimensionality reduction	PyTorch, TensorFlow, Keras	-
U-Net		Image segmentation, particularly in medical imaging	PyTorch, TensorFlow, Keras	-
CLIP		Cross-modal learning, image-text retrieval	OpenAI's CLIP, Hugging Face Transformers	Hugging Face Platform
DALL-E / Stable Diffusion		Text-to-image generation	OpenAI's DALL-E, Hugging Face Stable Diffusion	Hugging Face Platform
RAG		Open-domain question answering	Hugging Face Transformers, PyTorch	Hugging Face Platform
Neural Style Transfer		Artistic image generation, photo filters	PyTorch, TensorFlow, Keras	-
Deepfake Technology		Face swapping, video manipulation	PyTorch, TensorFlow, Keras	-
AGG		Enhanced text generation, multimodal content creation	PyTorch, TensorFlow, Hugging Face Transformers	Hugging Face Platform

GEN AI MODELS

Power Law Distribution of Gen AI

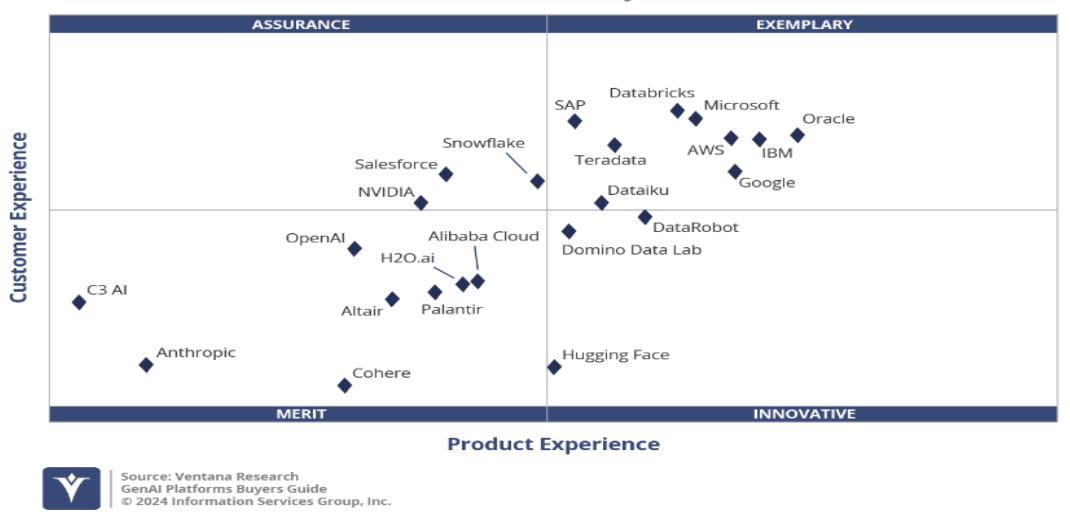


Generative foundation models

	OpenAI	Meta AI	Google	Other
Text	GPT-4	LLaMA	PaLM 2 Bard	Falcon Dolly Alpaca
Image	CLIP DALL-E 2	ImageBind SAM	Imagen	Stable Diffusion Midjourney
Audio	Whisper	Voicebox	MusicLM	

GenAI Platforms

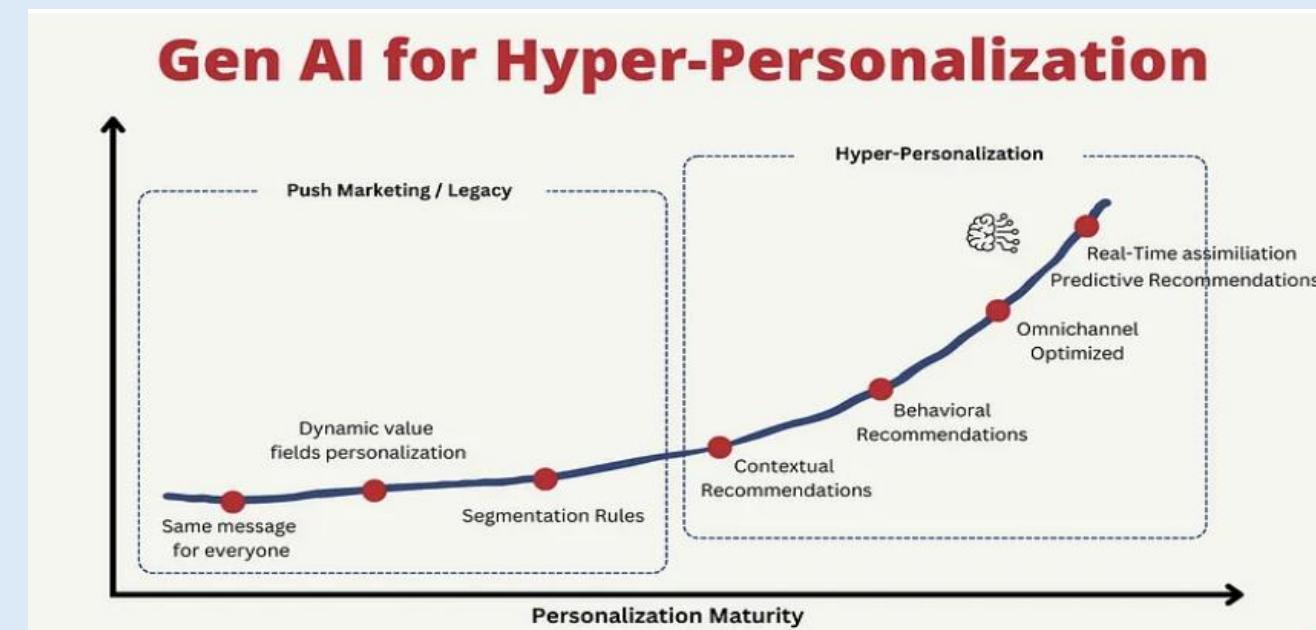
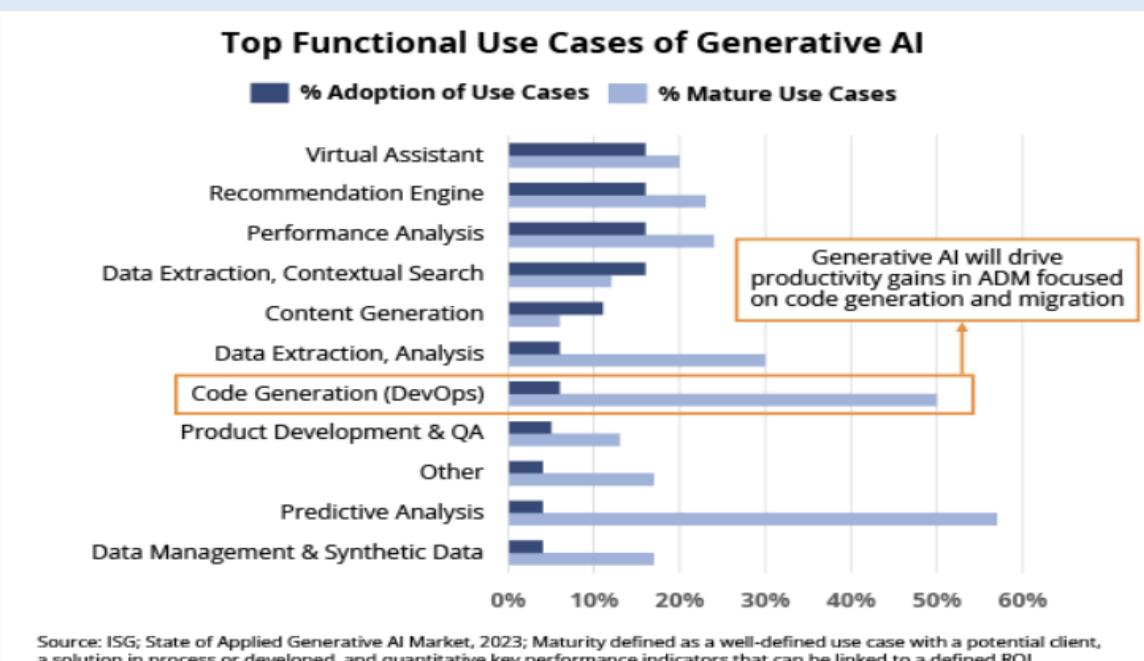
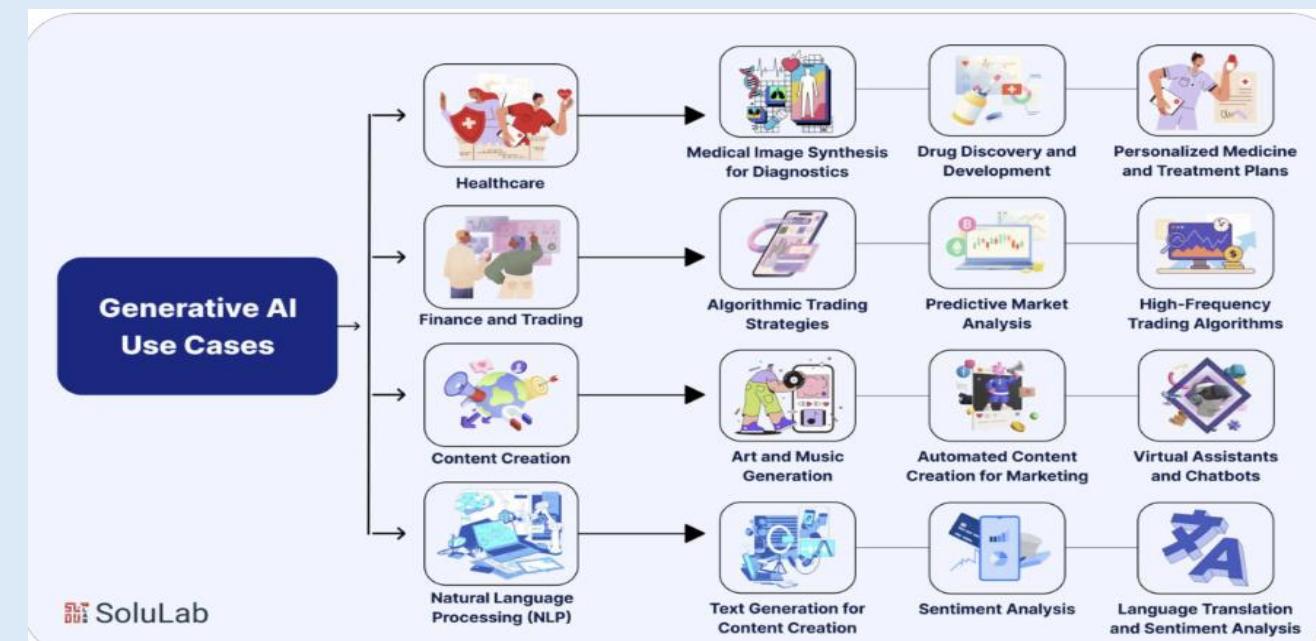
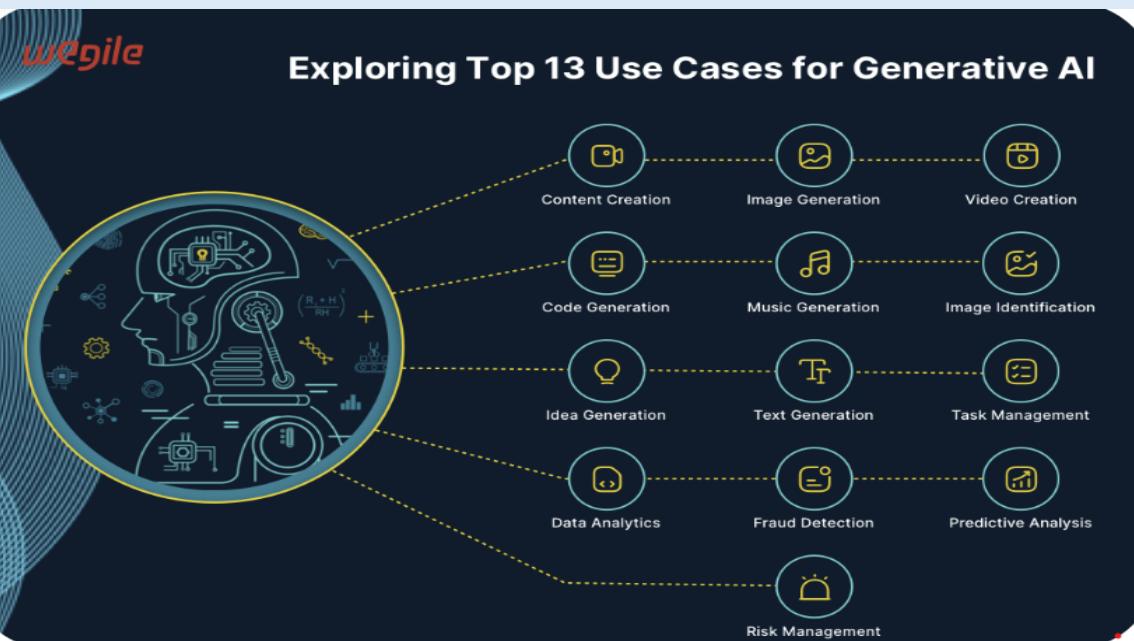
Ventana Research Buyers Guide



Model	Type	Use Cases
BERT	Free	Question Answering, Sentiment Analysis, NER
RoBERTa		Text Classification, Sentiment Analysis
DistilBERT		Sentiment Analysis, Text Classification
ALBERT		Sentence Classification, Sentiment Analysis
mBERT (Multilingual)		Multilingual Sentiment Analysis, Text Classification

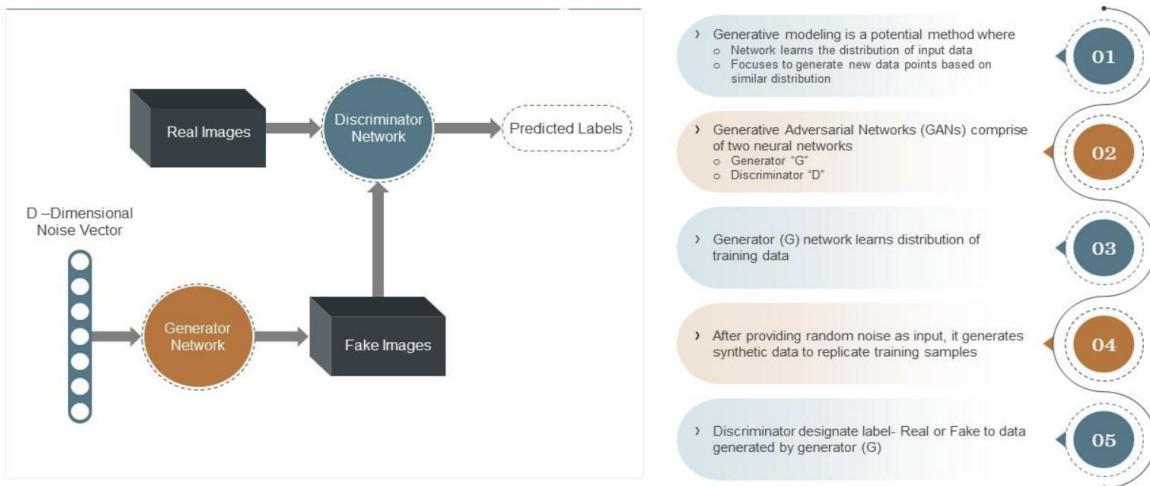
Model	Cost	Use Cases
GPT-4 (OpenAI)	Paid API (Subscription or per-token)	Text Generation, Chatbots, Summarization, QA
Claude 2 (Anthropic)		Text Generation, Summarization, Conversations
LLaMA 2 (Meta)		Language Modeling, Text Understanding, QA
BLOOM		Multilingual Text Generation, Summarization
T5 (Text-to-Text Transfer)		Summarization, Translation, Sentiment Analysis

GEN AI USE CASES



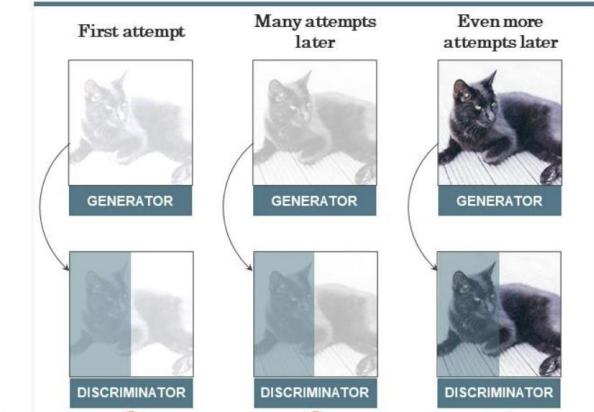
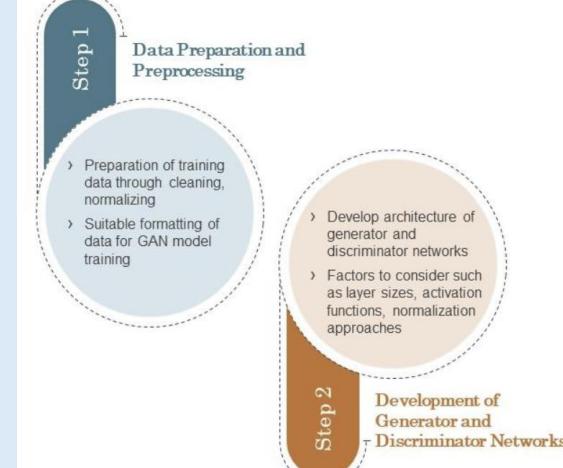
Architecture of generative adversarial networks (GANs)

This slide provides information regarding the framework of generative adversarial networks (GANs) which comprises of generator and discriminator. The network focuses on generating new data points based on similar distribution.



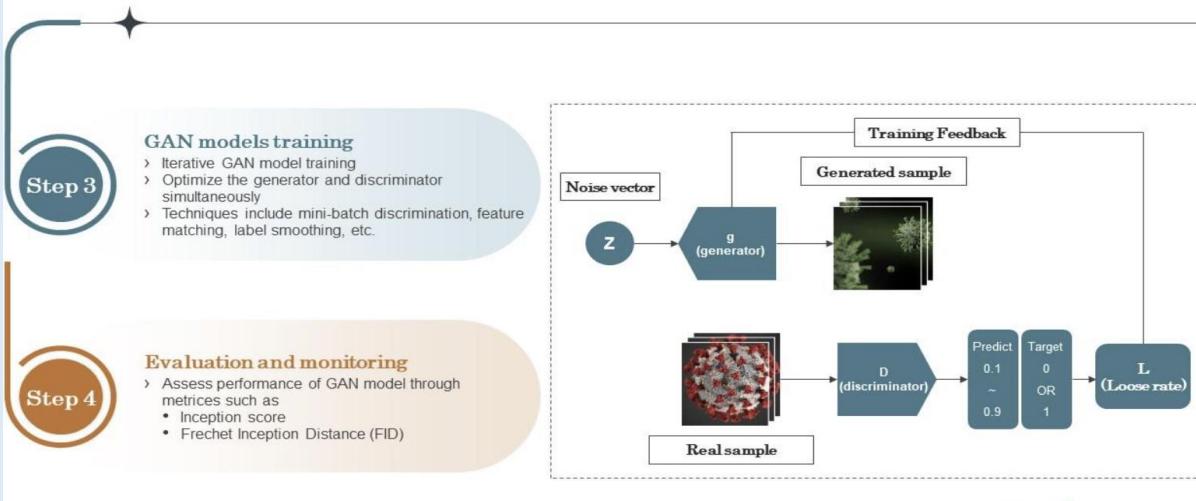
Essential steps for effective implementation of GAN models

This slide provides information regarding essential steps for the deployment of GAN models in terms of data preparation and preprocessing, and development of generator & discriminator networks.



Essential steps for effective implementation of GAN models cont.

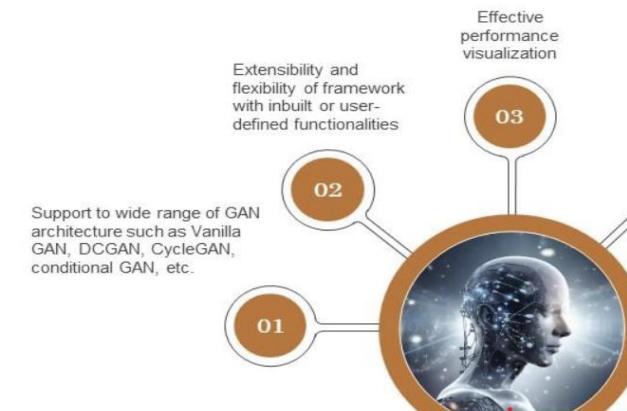
This slide provides information regarding essential steps for the deployment of GAN models in terms of data GAN models training along with evaluation and monitoring.



Popular GAN frameworks: PyTorch enabled Torch-GAN

This slide provides information regarding popular GAN framework as PyTorch enabled Torch-GAN suitable for building short and easy-to-manage codes while development of GANs, effective performance visualization, extensibility or flexibility of framework.

- 01: TorchGAN is PyTorch-based framework
- 02: Suitable for writing short and easy-to-manage codes for building GANs
- 03: Package comprises of several generative adversarial networks along with utilities essential for deployment
- 04: Torch GAN mimics GANs design through simple API that enables customization of components when required
- 05: Facilitates interaction among GAN components through highly flexible trainer that adopts user-defined model and losses



VAEs

Variational autoencoders (VAEs) as essential part of generative AI model

This slide provides information regarding variational autoencoders as a kind of generative model that builds upon conventional autoencoders. They are potential tools for learning compact and relevant latent representations of data while generating new data samples.

- › Variational autoencoder (VAE) is a kind of neural network that learns to regenerate its input and maps data to latent space
- › Such autoencoders have regularized training to limit overfitting and enable latent space have good properties to enable generative process



Architecture comprises of encoder and decoder which is trained to minimize reconstruction error among encoded-decoded data along with initial data

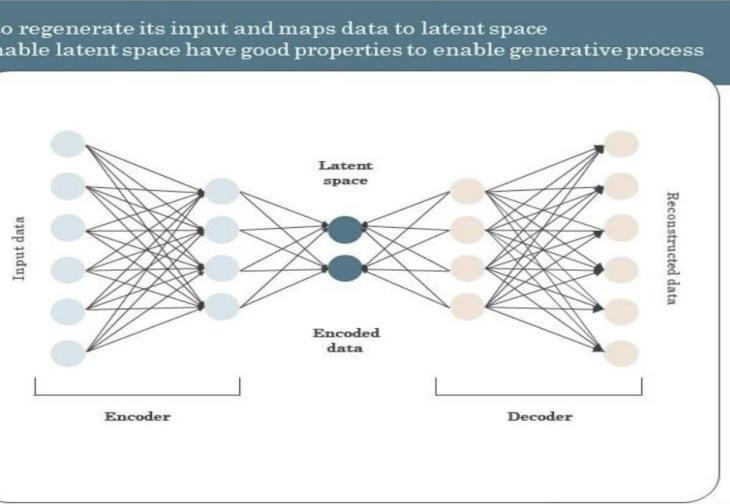
Input encoding is done as distribution instead of single point over latent space for introducing regularization

Regularization enables VAEs to generate new data samples that seamlessly interpolated among training data points

Enable VAEs as a potential tools for building new data samples similar to training data

Regularization also enables VAE to limit the decoder network from reconstructing the input data perfectly

Decoder network is forced to learn wider general data representation, that helps in enhancing competency of VAE to build new data samples

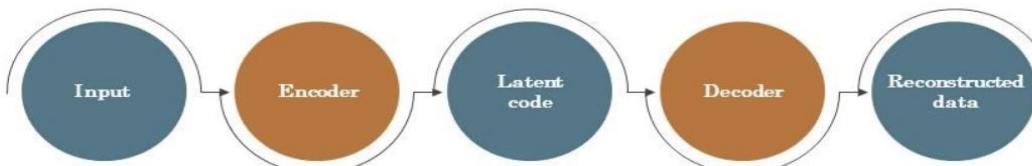


Training process of Variational Autoencoders (VAEs)

This slide provides information regarding training process of variational autoencoders which are suitable in generating new samples by learning from training dataset. The model learn from latent representation of data and sampling from latent space to create new samples.

VAEs are suitable to utilized as generative model, and are competent to generate new samples in context with training dataset

It is attained by learning a latent representation of data and sampling from the latent space to develop new samples



Process of generating new samples through VAE by training model

Step 1
Encode a given input sample into latent space

Step 2
Sample a new latent point from distribution defined by encoded sample

Step 3
Latent point is passed through decoder network to build output sample

Step 4
Reconstruction error is backpropagated across network as it learns to map latent points to data distribution on consistent basis

Step 5
Interpolating among latent points to generate new samples for smooth variations of input data



Pros of VAEs

Competent in generating new data samples by sampling from learned latent distribution and sample decoding

VAEs helps in reducing overfitting making them robust to handle complex datasets

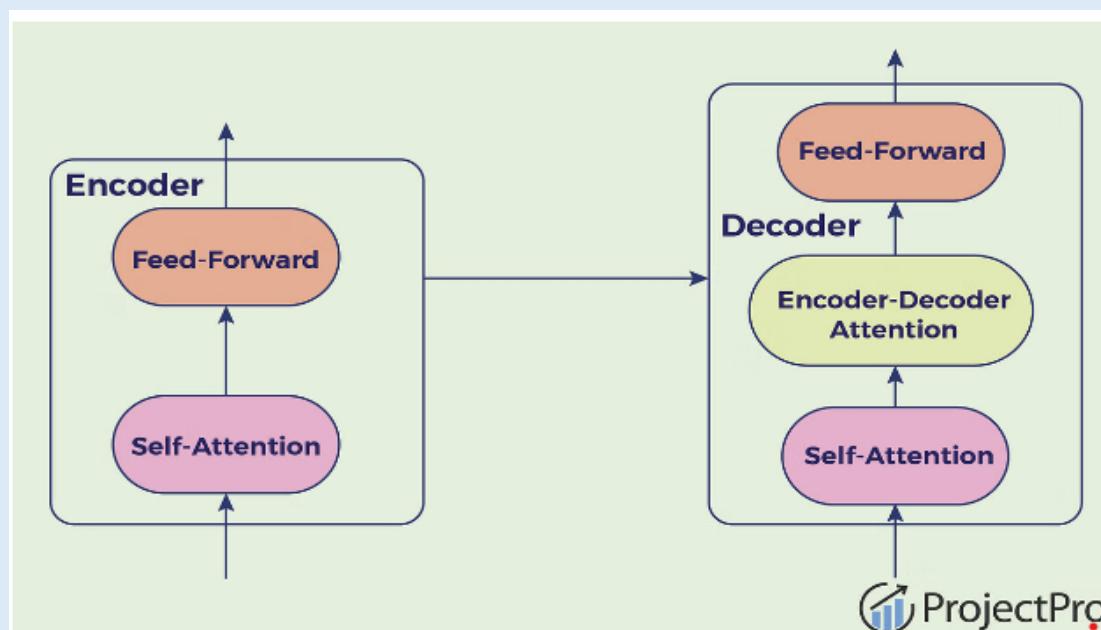
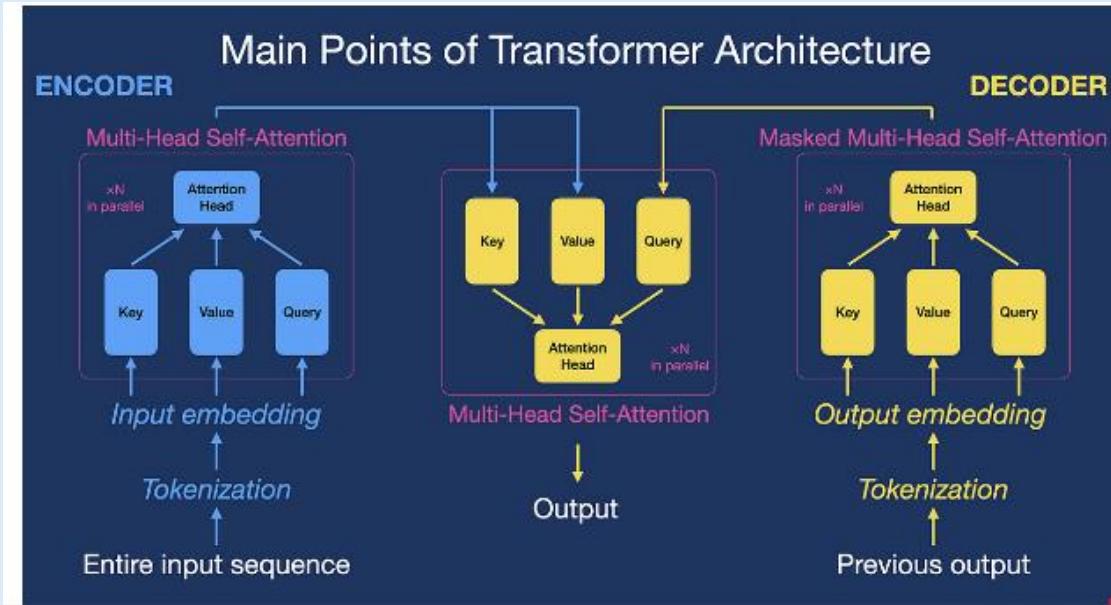


Cons of VAEs

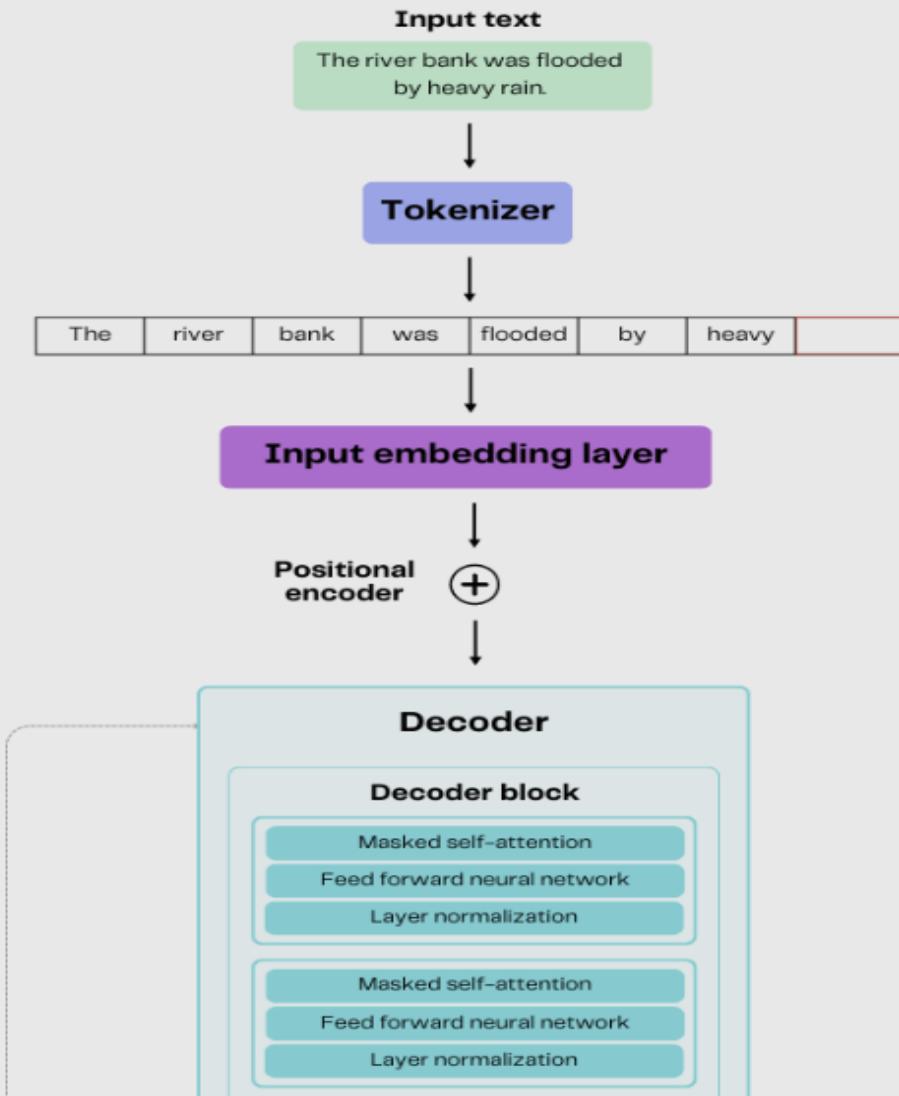
VAEs might generate blurry reconstructions

Complex training process

TRANSFORMER

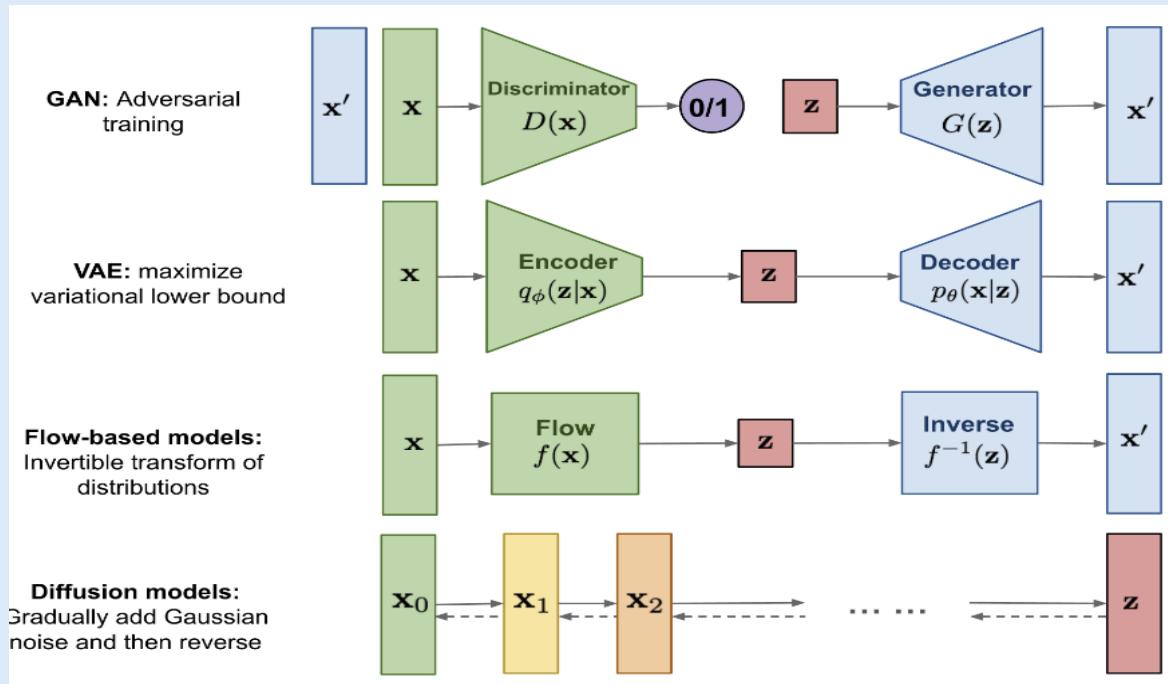
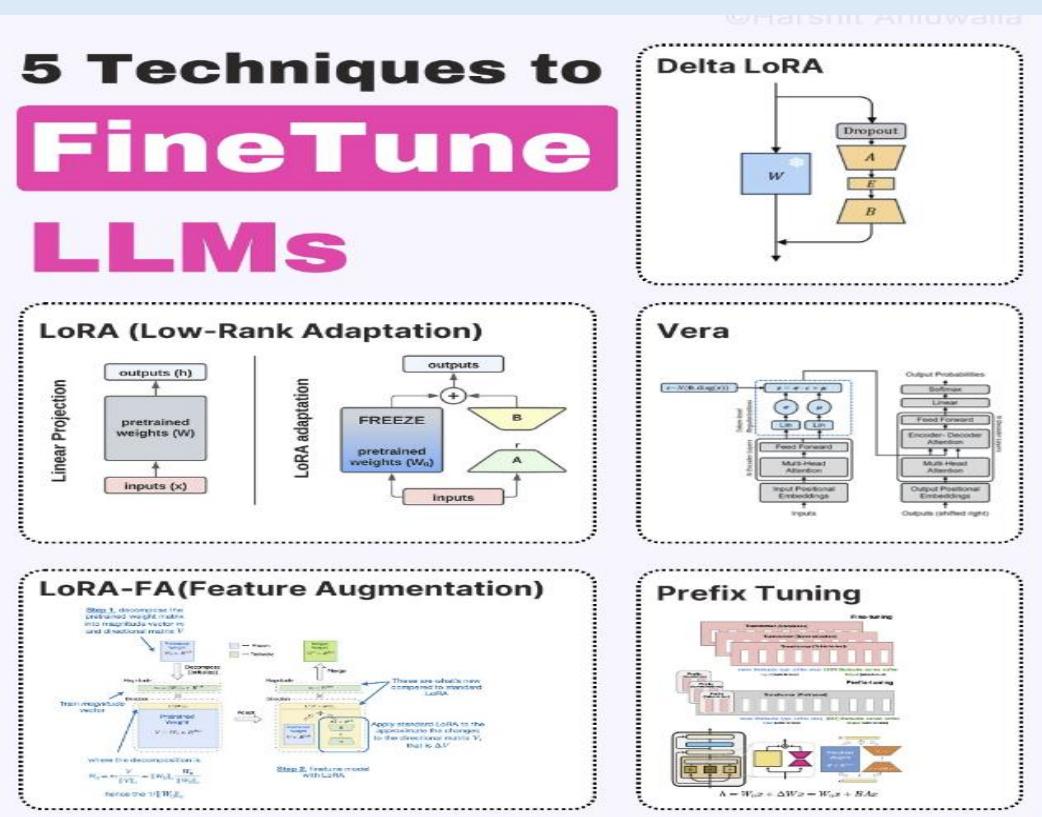
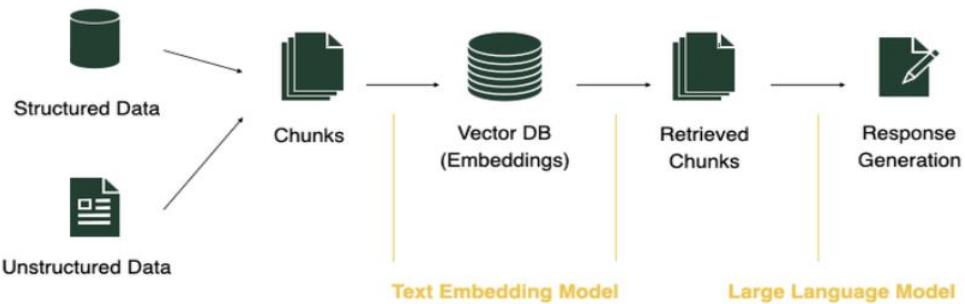


Pre-training transformers (decoder-only)



RAG, Diffusion model and Model Fine Tuning

Simple RAG



Pre-training diffusion models

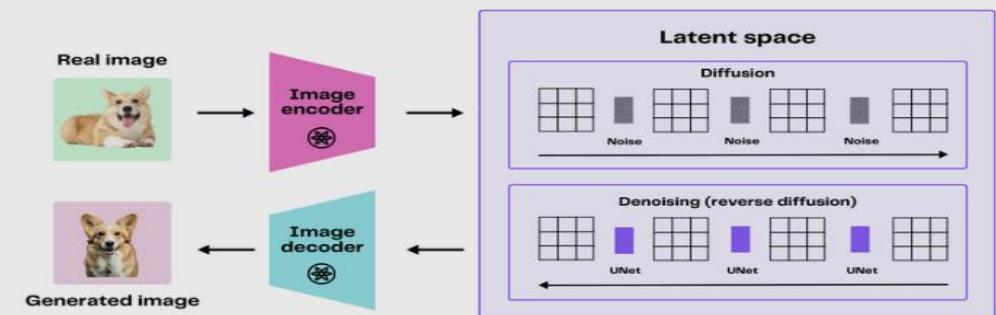




Figure 9: The LLM arms race with exponentially increasing parameter counts. (Credit: HuggingFace)

Parameters of transformer-based language models



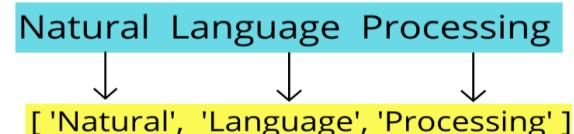
Natural Language Processing Techniques



Stemming vs Lemmatization



Tokenization



Increasingly longer prompt

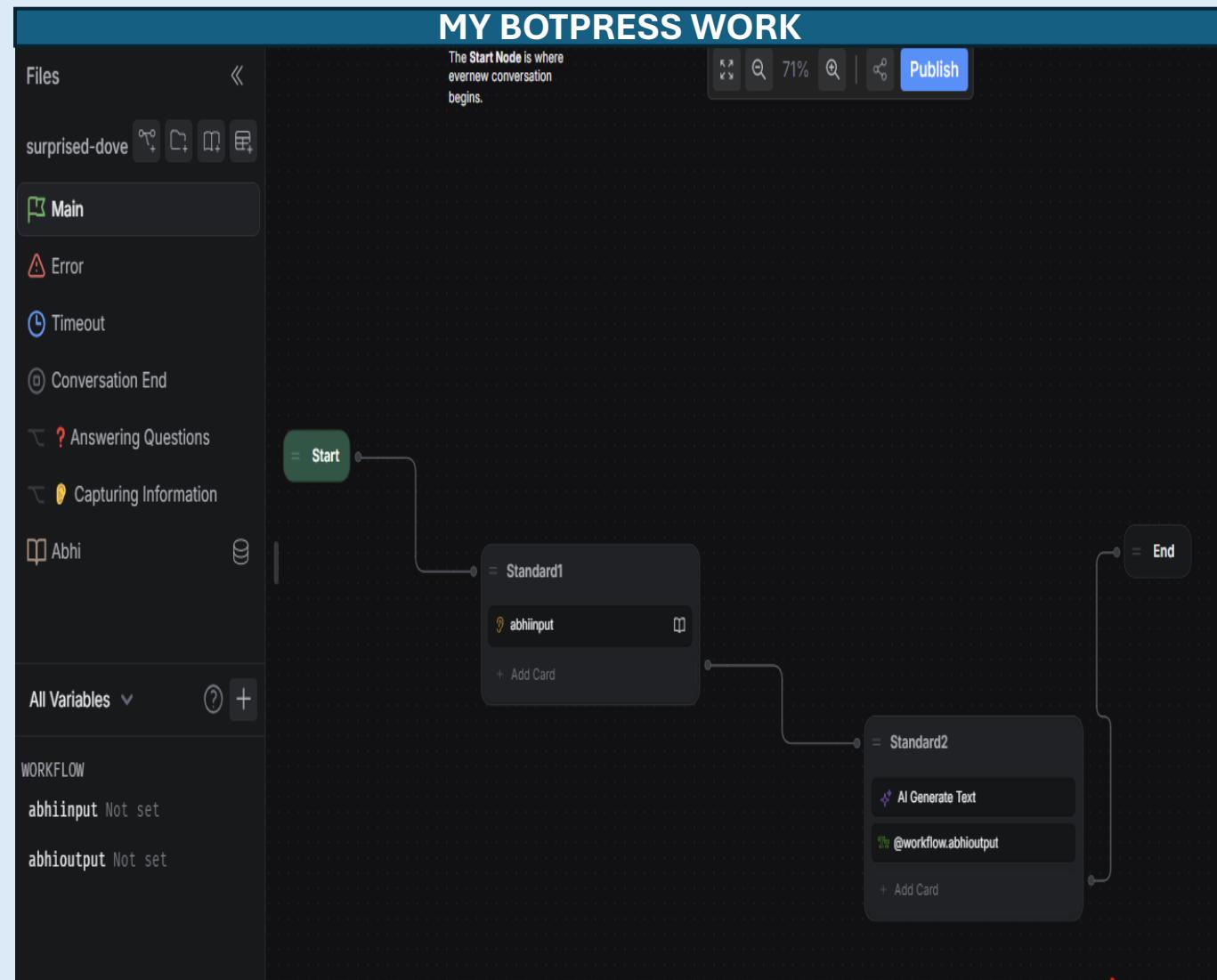
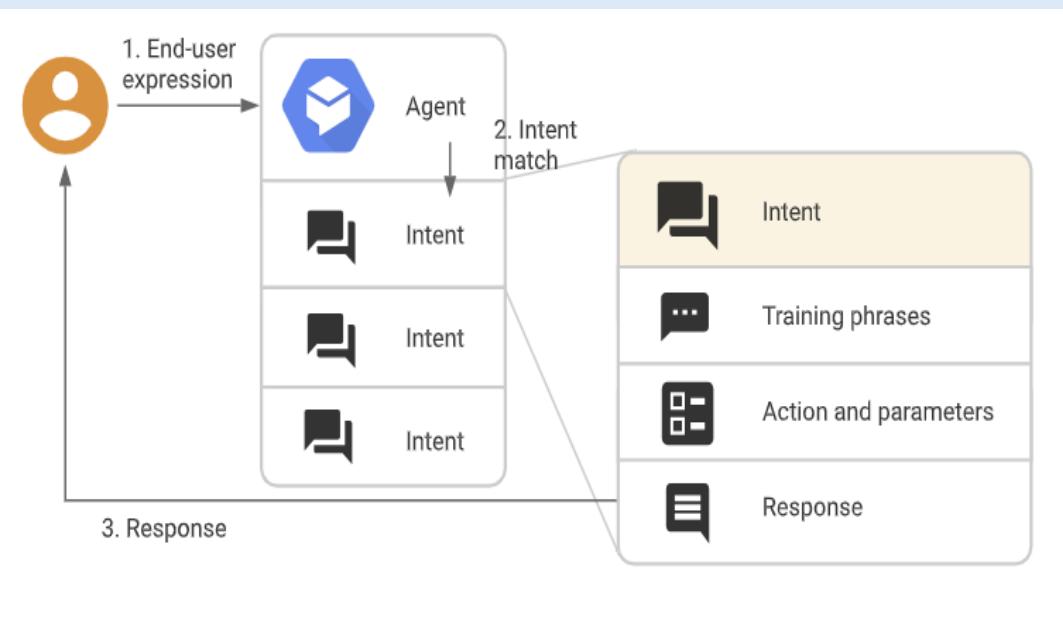
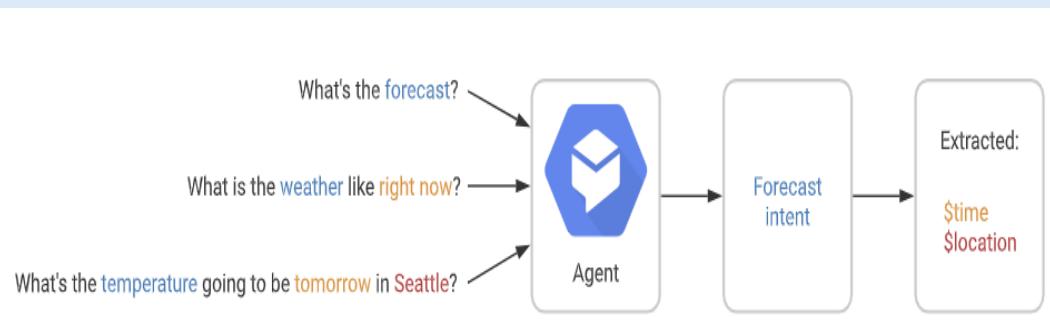
- ─ In-Context Learning, CoT
- ─ Retrieval Augmented Generation
- ─ Copilots, Tools, Agents,
- ─ More than 20k tokens

Natural Language



- Challenges**
- ─ High latency, including multi-call;
 - ─ Limited context windows;
 - ─ Forgot context;
 - ─ Hefty cost;
 - ─ Performance drop, like lost in the middle;

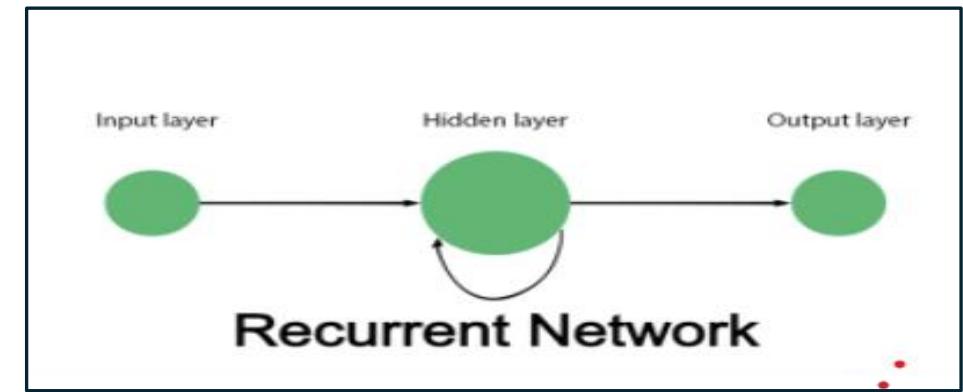
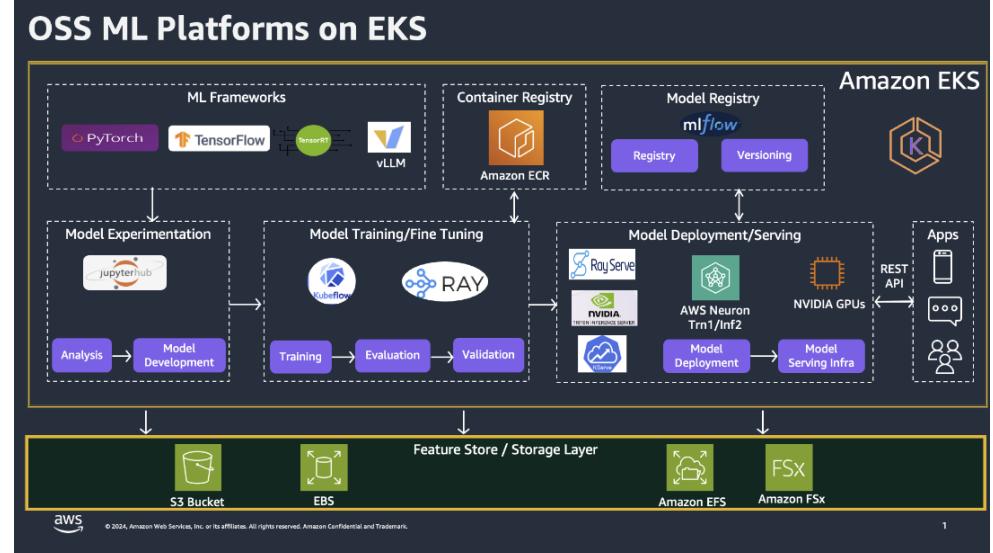
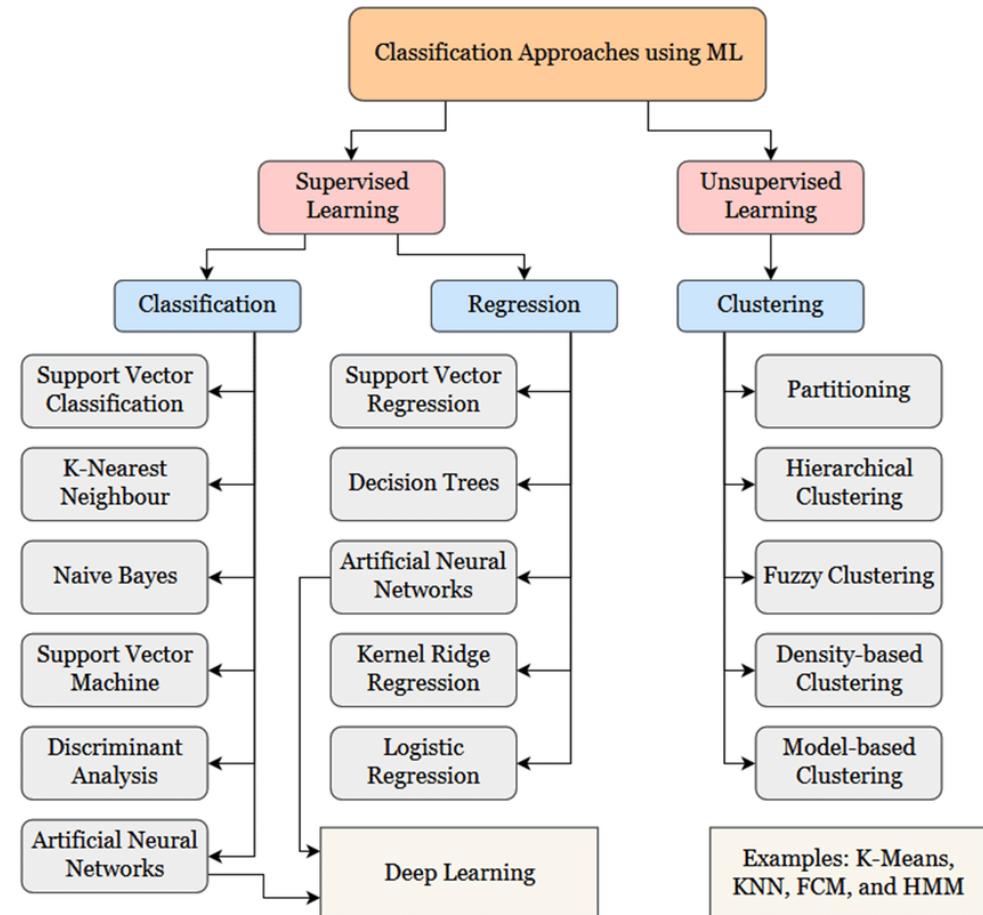
DIALOGFLOW/ BOTPRESS



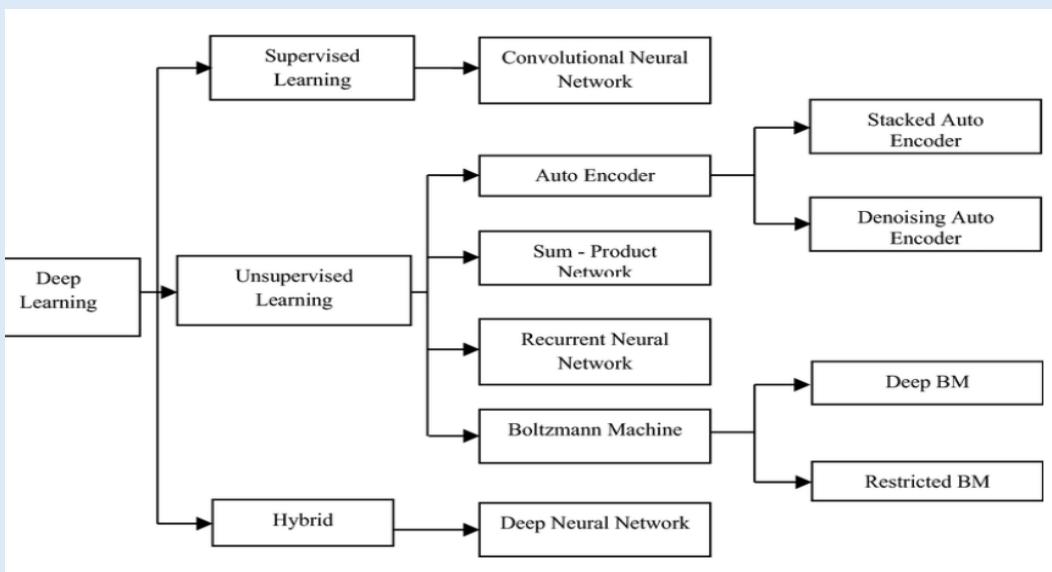
THANKYOU

Github Repository link --<https://github.com/Abhishekrai129>

APPENDIX



APPENDIX

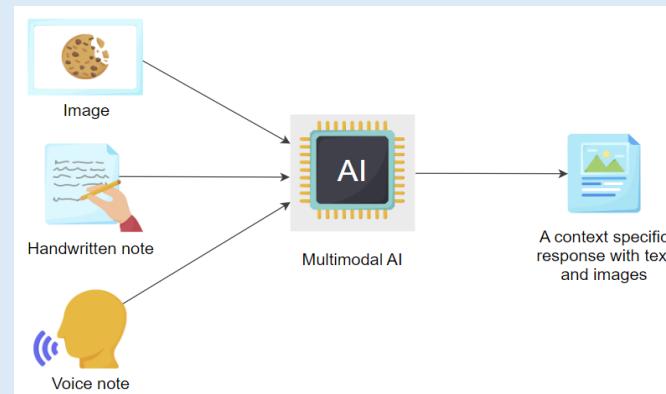
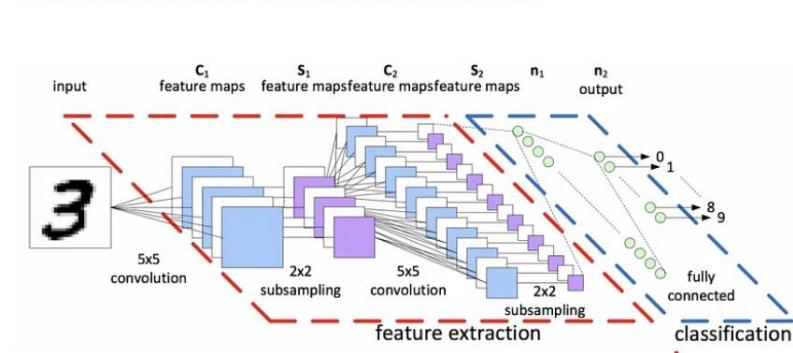


Building Blocks: 7 AI Primitives

Primitive	Closed-source model example	Open-source model example	Application example
① Language models	ANTHROPIC Claude	G BERT	Lindy
② Image generation models	DALL-E 2	stability.ai	Typeface
③ Video creation models	runway	deforum-stable-diffusion	Hour One
④ Video indexing models	Twelve Labs	MERLOT	VOLT
⑤ Speech synthesis models	ElevenLabs	tortoise-tts-v2	synthesia
⑥ Speech comprehension models	DEEPGRAM	OpenAI Whisper	Speak
⑦ Vector databases	Pinecone	Weaviate	mem

MENLO VENTURES

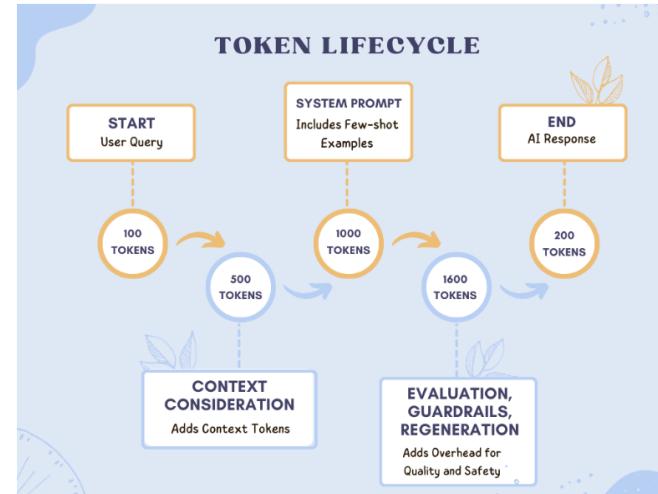
Convolutional Neural Network Architecture



Feature	ANN	RNN	CNN
Structure	Feedforward	Recurrent (loops)	Convolutional layers
Input Type	Structured data	Sequential data	Grid-like data (images)
Memory	Stateless	Maintains state (memory)	Stateless
Use Cases	General tasks	NLP, time-series	Image processing
Training Complexity	Moderate	Higher (vanishing gradients)	Moderate

APPENDIX

GEN AI MODELS



Model	Type	Use Cases
BERT	Free	Question Answering, Sentiment Analysis, NER
RoBERTa		Text Classification, Sentiment Analysis
DistilBERT		Sentiment Analysis, Text Classification
ALBERT		Sentence Classification, Sentiment Analysis
mBERT (Multilingual)		Multilingual Sentiment Analysis, Text Classification
Model	Cost	Use Cases
GPT-4 (OpenAI)	Paid API (Subscription or per-token)	Text Generation, Chatbots, Summarization, QA
Claude 2 (Anthropic)		Text Generation, Summarization, Conversations
LLaMA 2 (Meta)		Language Modeling, Text Understanding, QA
BLOOM		Multilingual Text Generation, Summarization
T5 (Text-to-Text Transfer)		Summarization, Translation, Sentiment Analysis

Category	Top 5 Models	Frameworks
1. NLP & LLM (Natural Language Processing)	1. GPT-3 / GPT-4 (OpenAI) 2. BERT (Google) 3. LLaMA (Meta) 4. T5 (Google) 5. RoBERTa (Facebook)	1. Transformers (Hugging Face) 2. TensorFlow / Keras 3. PyTorch 4. Fairseq 5. AllenNLP
2. Computer Vision	1. YOLO 2. ResNet 3. EfficientNet 4. Mask R-CNN 5. Vision Transformer (ViT)	1. OpenCV 2. TensorFlow / Keras 3. PyTorch 4. Fastai 5. Detectron2
3. Audio and Speech Processing	1. DeepSpeech 2. Tacotron 2 3. WaveNet 4. VITS 5. HuBERT	1. Kaldi 2. TensorFlow / Keras 3. PyTorch 4. ESPnet 5. Librosa
4. Generative Models (GANs, VAEs, etc.)	1. StyleGAN 2. CycleGAN 3. VAE 4. DALL-E 5. Stable Diffusion	1. TensorFlow / Keras 2. PyTorch 3. Fastai 4. NVIDIA StyleGAN2 5. CLIP
5. Reinforcement Learning & Robotics	1. DQN 2. PPO 3. A3C 4. TD3 5. OpenAI Five	1. OpenAI Gym 2. Ray RLLib 3. Stable Baselines3 4. TF-Agents 5. PyTorch Lightning

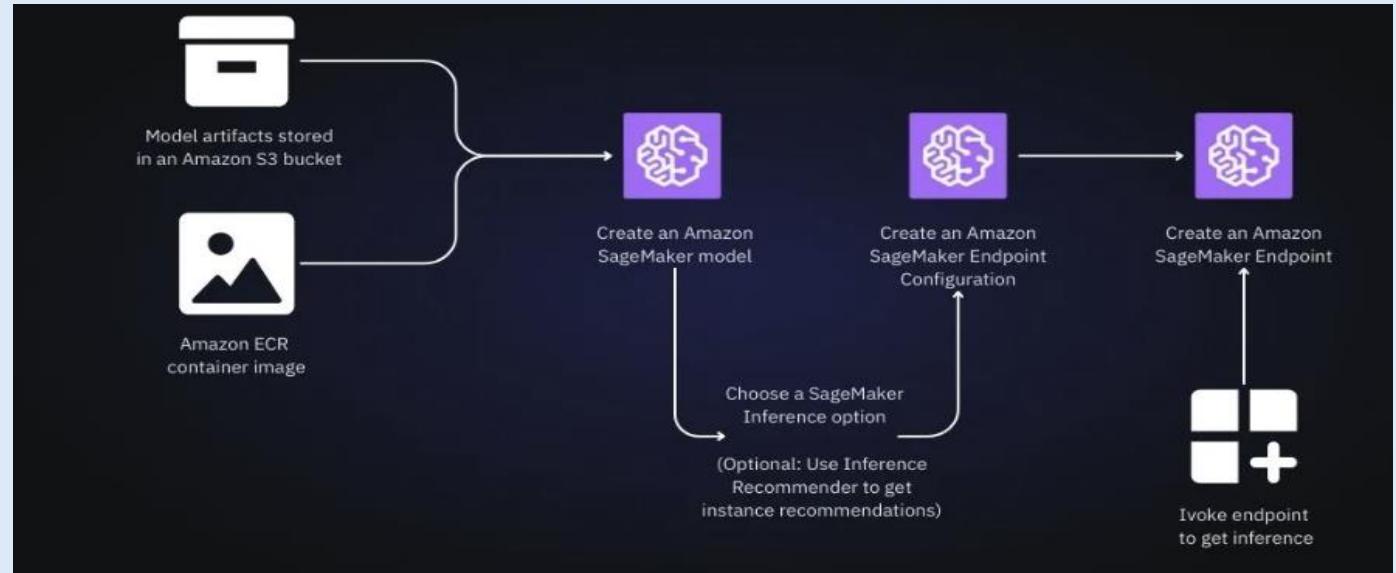
APPENDIX

Company	Model	Type	Price
OpenAI	GPT 4.0 8k context	Input	\$0.03/1k tokens
		Output	\$0.06/1k tokens
	GPT 4.0 32k context	Input	\$0.06/1K tokens
		Output	\$0.12/1K tokens
	GPT-3.5 Turbo 4k context	Input	\$0.0015/1K tokens
		Output	\$0.002/1K tokens
Google	GPT-3.5 Turbo 16k context	Input	\$0.003/1K tokens
		Output	\$0.004/1K tokens
	PaLM 2 for Text	Input	\$0.0005/1K characters
		Output	\$0.0005/1K characters
PaLM 2 for Chat	Input	\$0.0005/1K characters	

The top datasets are all based on images with a few text-based datasets such as Wikipedia and PubMed.

Table 1 Top 20 datasets mentioned in the open access subset of the GenAI corpus

Dataset name	Citing documents	Total mentions	Main modality
1 ImageNet	2,741	6,823	image
2 MNIST	2,533	9,292	image
3 CIFAR-10	2,160	7,744	image
4 CelebA	1,705	5,713	image
5 COCO	1,141	3,390	image
6 Wikipedia	662	2,599	text
7 FFHQ dataset	596	1,983	image
8 FASHIONMNIST	520	1,375	image
9 CelebAHQ	474	1,081	image
10 SVHN	398	1,414	image



Strategies for Cost Management in GEN AI Implementation

