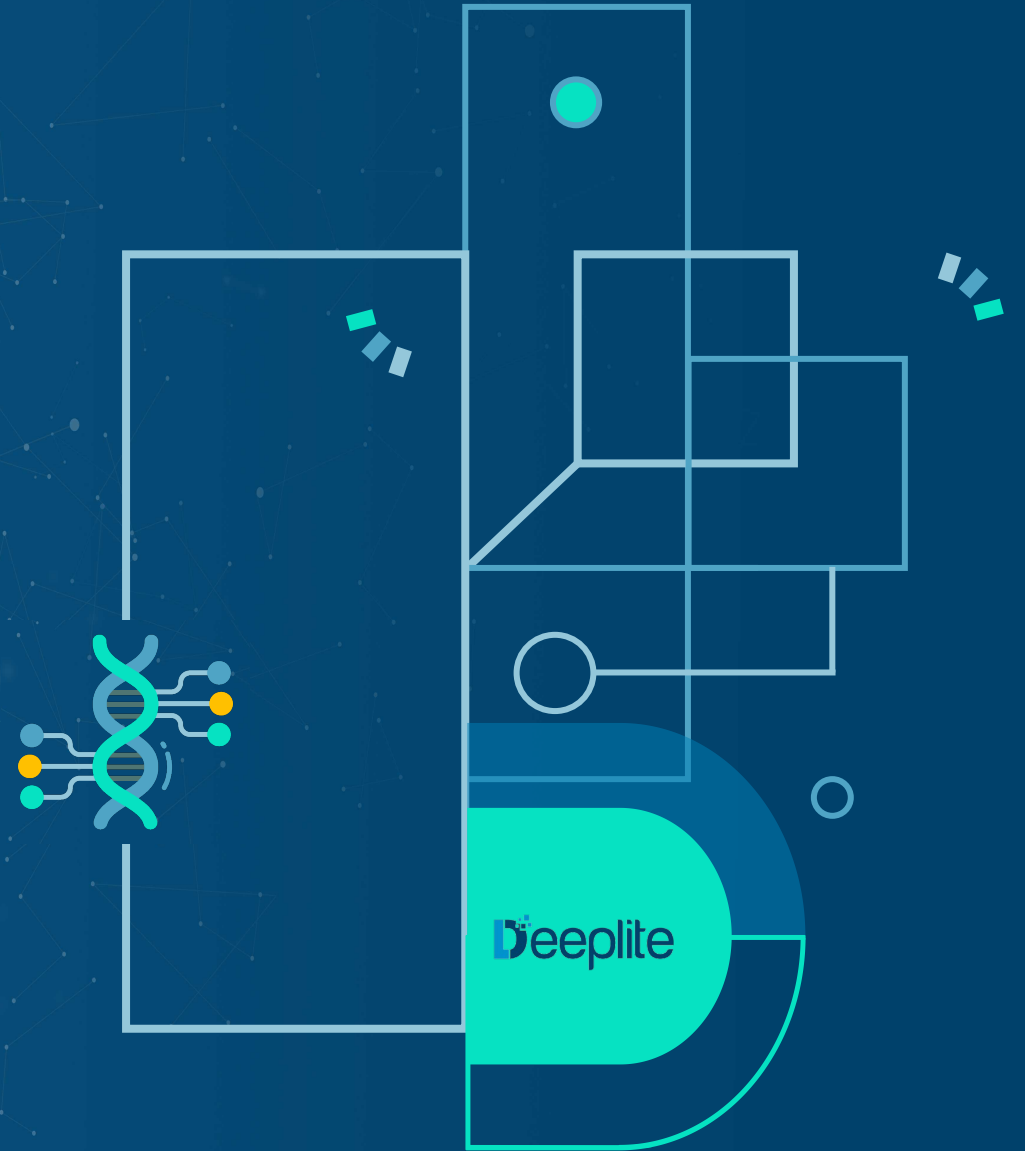# Deeplite

Optimizing AI for the Edge

## Introduction to Artificial Intelligence and It's Impact on Business

**April 24, 2024**

Confidential

**Nick Romano**
www.linkedin.com/in/nickromanoprofile/
nick@deeplite.ai

Cofounder and CEO of Deeplite, Nick is a serial entrepreneur and 3-time founder. He has Bachelor of Engineering and Management degree in Mechanical Engineering from McMaster University, Canada. Recently honored by McMaster's Engineering Faculty as being one of their Top 150 Distinguished Alumni for the role they've played in shaping Canada and the world.
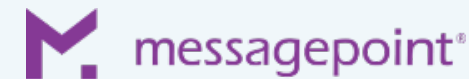
# CEO | Director | Advisor | Private Equity | SaaS | AI

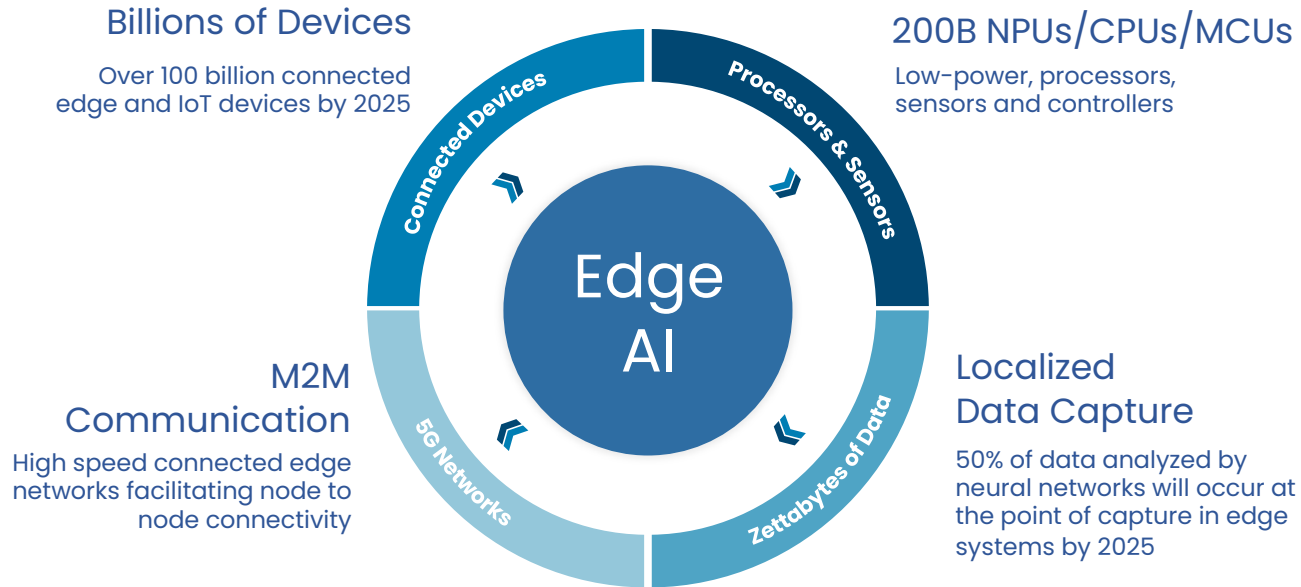| | | |
|---|---|---|
| Deeplite | Canadian Chamber of Commerce | McMaster University |
| 2019 | 2023 | 2024 |
| prinova | messagepoint | |
| 1998 | 2008 | |

# AI at the Edge
Opportunity and challenges

**Deeplite**

AI-Driven Optimization to make Deep Neural Networks faster, smaller and energy-efficient. **Edge will be bigger than cloud**

**Billions of Devices**

Over 100 billion connected edge and IoT devices by 2025

**200B NPUs/CPUs/MCUs**

Low-power, processors, sensors and controllers

## Edge AI

*Connected Devices*
*Processors & Sensors*
*5G Networks*
*Zettabytes of Data*

**M2M Communication**

High speed connected edge networks facilitating node to node connectivity

**Localized Data Capture**

50% of data analyzed by neural networks will occur at the point of capture in edge systems by 2025

**AI becomes untethered, decentralized and everywhere**

## Investors

PJC

INNOSPARK
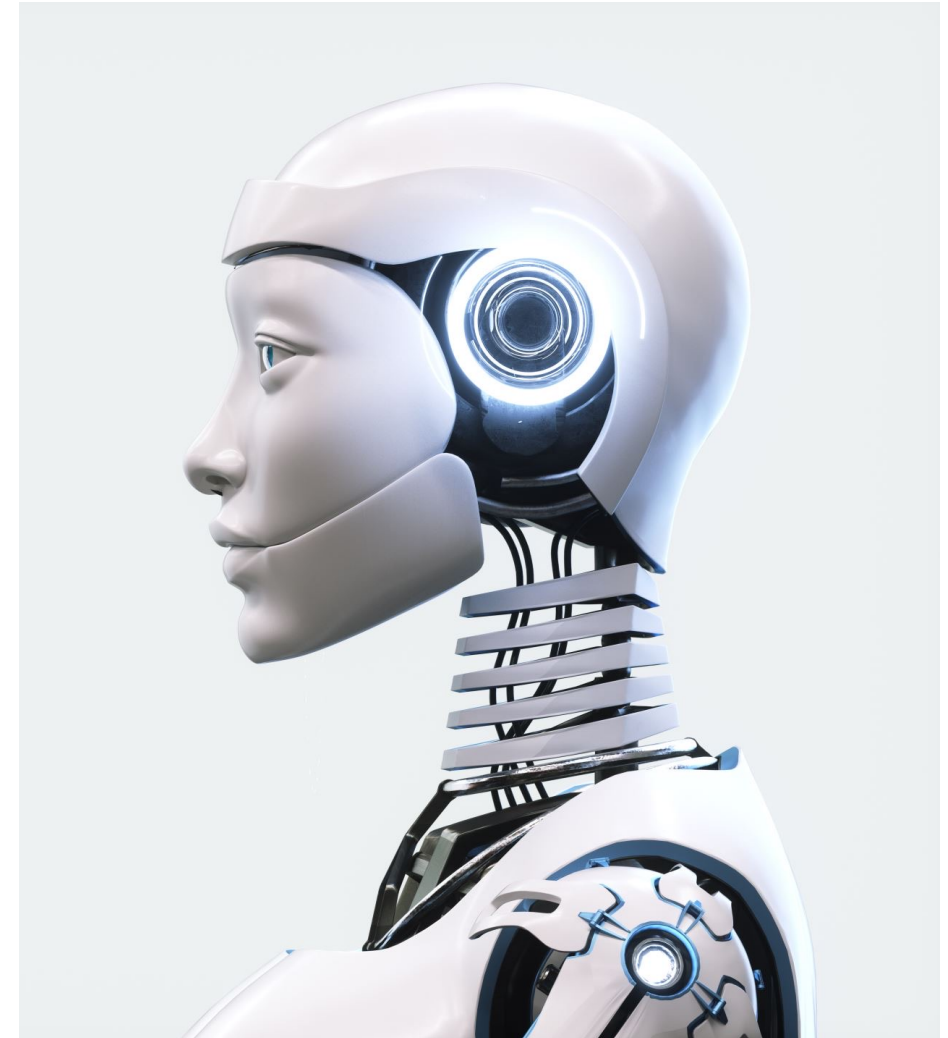
DIFFERENTIAL VENTURES

bdc

Desjardins Capital

SOMEL

TANDEMLAUNCH

# What is AI?
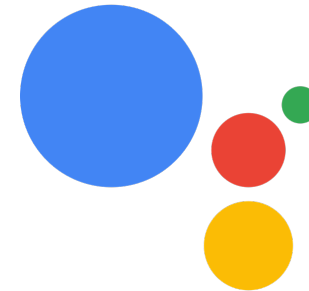
- AI is a branch of computer science that aims to create intelligent machines.

- AI is about creating systems that can perform tasks that would require human intelligence.

# AI in Everyday Life

- Personal Assistants
- Recommendation Systems
- Navigation and Travel
- Health and Fitness
- Customs
- Many more...

# Types of AI

**Deeplite**

## Narrow AI

- "Weak AI"
- Designed to perform a specific task such as "Person Detection"
- Focused with limited constraints
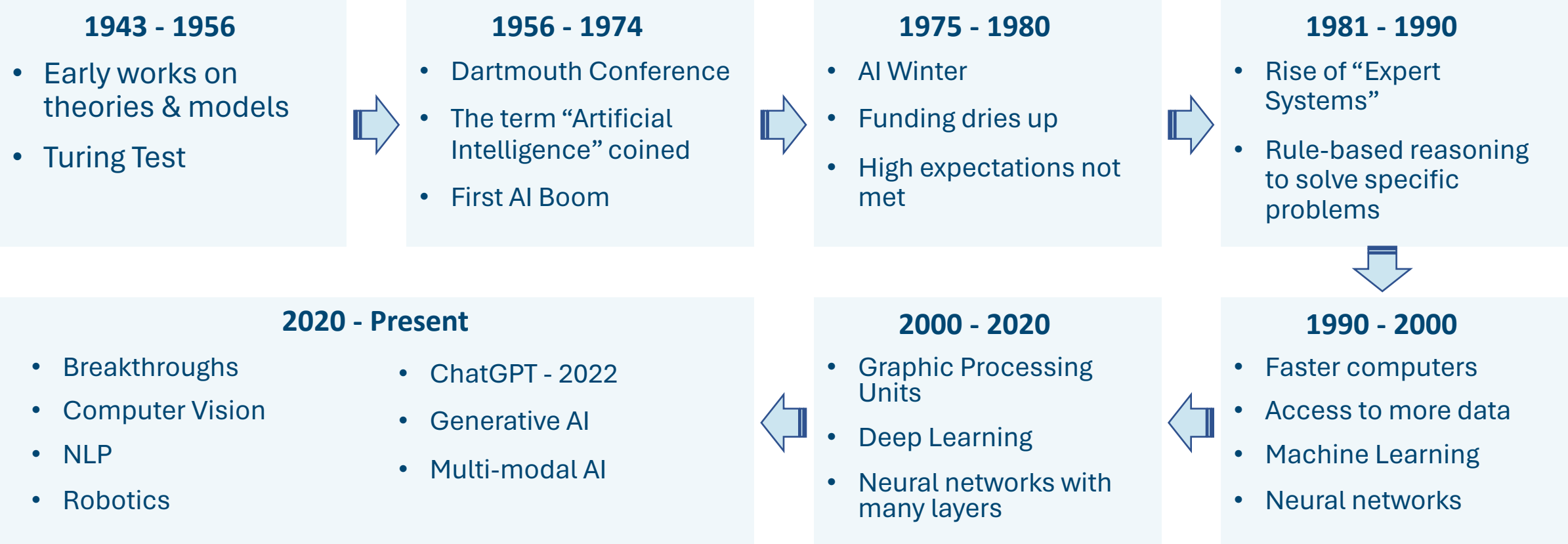- Lacks consciousness, genuine understanding and self-awareness
- This is AI today

## General AI

- "Strong AI" or "AGI"
- Can perform like a human
- Can understand, learn, adapt
- Implement knowledge from one domain to another
- Understands context and make judgements based on that
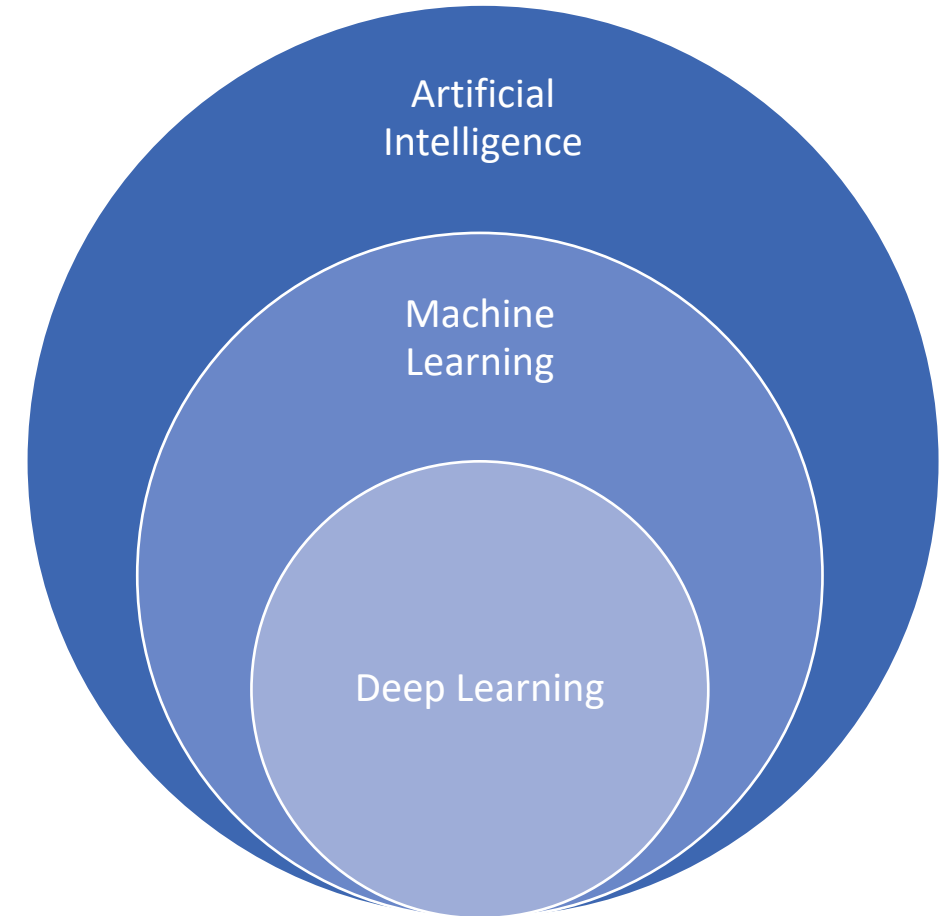- Human level cognitive function
- Does not exist "yet"

## Superintelligent AI

- Surpasses human intelligence
- Self-improvement
- Could solve complex problems that human couldn't
  - Curing diseases
  - Predicting stock market
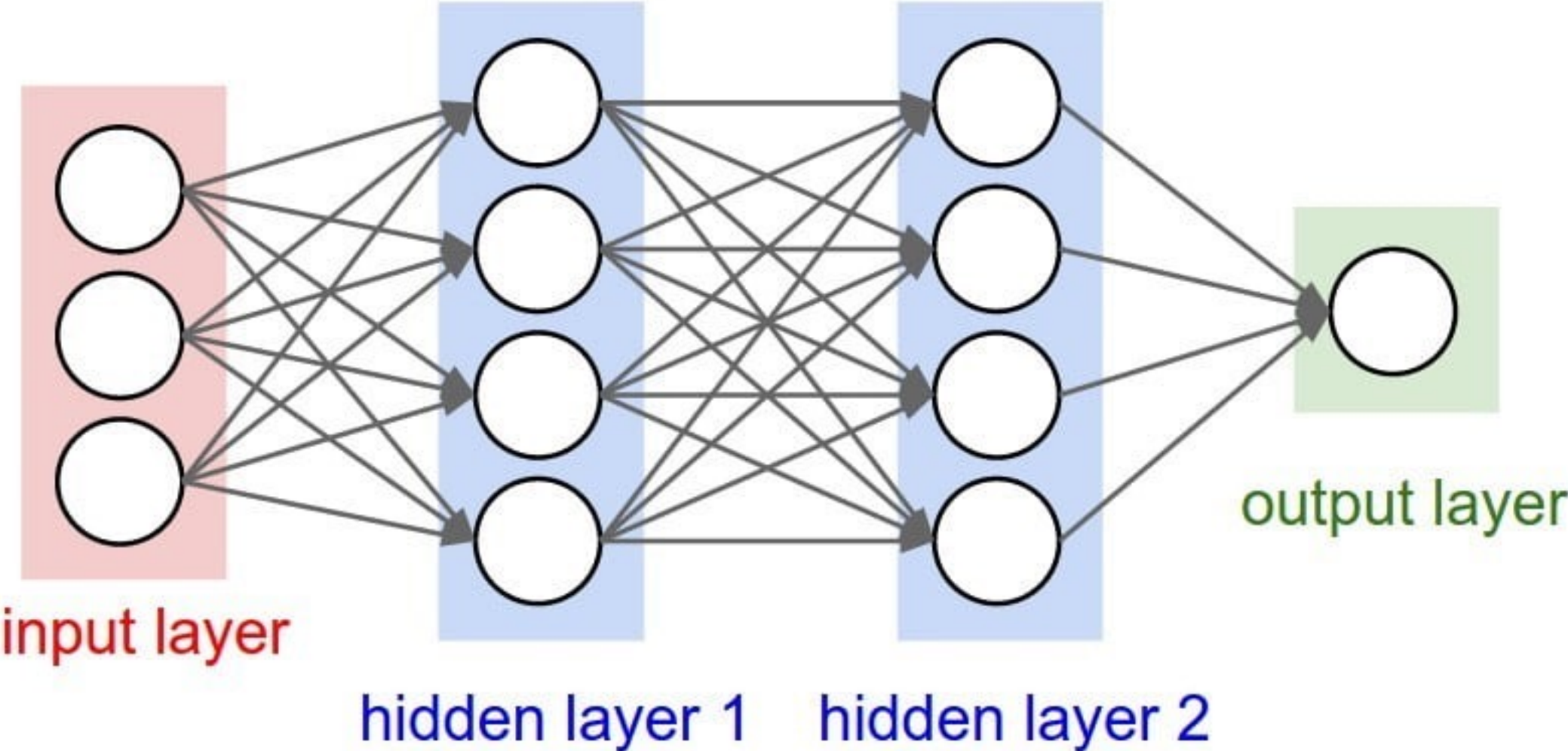- Major risk & ethical considerations
- Hypothetical

# History of AI

**Deeplite**

## 1943 - 1956

- Early works on theories & models
- Turing Test

## 1956 - 1974

- Dartmouth Conference
- The term "Artificial Intelligence" coined
- First AI Boom

## 1975 - 1980

- AI Winter
- Funding dries up
- High expectations not met

## 1981 - 1990

- Rise of "Expert Systems"
- Rule-based reasoning to solve specific problems

## 2020 - Present

- Breakthroughs
- Computer Vision
- NLP
- Robotics
- ChatGPT - 2022
- Generative AI
- Multi-modal AI

## 2000 - 2020

- Graphic Processing Units
- Deep Learning
- Neural networks with many layers

## 1990 - 2000

- Faster computers
- Access to more data
- Machine Learning
- Neural networks

# How does AI Work?

- Machine Learning: AI learns from data to make predictions or decisions.

- Deep Learning: A subset of machine learning that uses neural networks with many layers. Requires LOTS of data to learn

- Reinforcement Learning: AI learns by trial and error to achieve a clear objective.

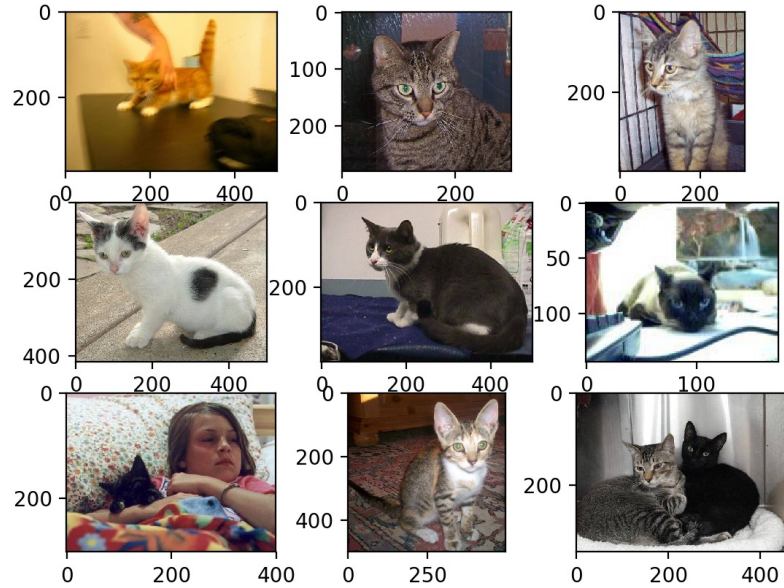- Imitation Learning: AI learns from a catalog of "demonstrations"

Artificial Intelligence

Machine Learning

Deep Learning

# Neural Network

Based on the operation of neurons in the human brain

input layer

hidden layer 1     hidden layer 2

output layer

# Training AI

Supervised and unsupervised learning

More data + bigger models = Better accuracy

Better accuracy = Better decisions

## AI is Resource Intensive

Human brain requires ~20 watts of energy

An AI100 NVIDIA GPU can consume up to 300 watt-hours (Wh) of energy

(7.5MW per hour)

# Training ChatGPT-3
## (175B parameters)

- 1,023 GPUs at ~34 days
- 250MW consumed
- Estimated 500+ metric tons of $CO_2$ released[1]

**ChatGPT-4**

- Is said to have 2T parameters!
- ~4-5,000 GPUs were used to train

1. Stanford Institute for Human-Centered AI

**Nvidia's high-end chips will consume the same amount of energy as a small nation in 2024**

Annual electricity consumption

| | | | | |
|---|---|---|---|---|
| 13,797 GWh | 13,000 GWh | 13,000 GWh | 11,000 GWh | 10,632 GWh |
| NVIDIA H100 Chips | Georgia | Guatemala | Costa Rica | 1M US homes |

Source: CB Insights estimates, assuming 3M H100s are operational in 2024 (based on 2024 sales estimate from Nvidia and Q3'23 sales estimate from Omdia Research) at 75% of max power; electricity consumption of countries and US homes based on 2021 EIA data

5 / CBINSIGHTS

# The Cost of Training AI

# AI Technologies

## Mastering the 5 senses

- Natural Language Processing: Allows machines to understand and respond to human language.

- Audio: Enables machines to hear and understand the content of sound.

- Computer Vision: Enables machines to see and understand the content of digital images or videos.

- Robotics: Design, construction, and use of robots to perform tasks done traditionally by human beings.

- Autonomous Vehicles: Use of AI in self-driving cars.

# AI that "generates" content based on a request or "prompt"
# • ChatGPT – Prompt to Text (written or voice exchange)

# The Answer Era



**Card Catalogue Era**



**Search or "Blue Links" Era**



**Answer Era**

Confidential

Deeplite

# AI that "generates" content based on a request or "prompt"

- ChatGPT – Prompt to Text
- Dall-e – Prompt to Image



Prompt: "A 3D render of a coffee mug placed on a window sill during a stormy day. The storm outside the window is reflected in the coffee, with miniature lightning bolts and turbulent waves seen inside the mug. The room is dimly lit, adding to the dramatic atmosphere." | Image: DALL-E 3 prompted by OpenAI

Deeplite

AI that "generates" content based on a request or "prompt"
- ChatGPT – Prompt to Text
- Dall-e – Prompt to Image
- Sora – Prompt to Video

Historical footage of California during the gold rush

**Deeplite**

AI that can understand and create content across different modalities like text, images, audio etc.

Combines perception and generative AI

One step closer to AGI!

# AI in Business

## Overview

- AI is transforming various business functions: marketing, HR, finance, legal etc.

- AI can automate repetitive tasks, freeing up time for more strategic work.

- AI can provide insights from large amounts of data to inform decision-making.

- AI can improve customer service through chatbots and personalized recommendations.

# Amazon

## How Amazon uses AI for Product Recommendations

- AI algorithms identify patterns in customer data

- AI generates personalized product recommendations.

- Use AI to monetize customer data.

- Algorithm generates product pitches
  - Help customers find what they are looking for
  - Nudge them towards buying more.

- Amazon's AI analyzes user behavior by tracking various data points
  - Products viewed,
  - Time spent on each page,
  - The frequency of purchases.

# Starbucks

How Starbucks uses AI to Personalize Customer Experiences

- Personalized customer experiences by leveraging big data and predictive analytics.

- The **Deep Brew** initiative, a platform that enables the seamless deployment of AI models across the company's expansive operations.

- AI algorithms built into the Starbucks app

- Analyzes customer preferences and behavior

- Provide personalized drink recommendations and offers.

- AI-powered chatbots provide customers with personalized recommendations and assistance.



STARBUCKS® REWARDS

FESTIVE THURS★YAYS

HALF OFF FOR THE HOLIDAYS

Enjoy half off a drink, every Thursday in December.
For Starbucks® Rewards members.

THURSDAYS | ALL DECEMBER | 12–6 PM

# American Express

How Amex uses AI for Fraud Detection

- AI powered fraud detection models

- Monitors each transaction in real-time

- Generates a fraud decision in milliseconds every single time an American Express card is used globally.

- American Express uses advanced AI models to detect anomalous patterns in transactions.

- The real-time fraud detection system improves accuracy, and better protect customers and merchants.

# Google

## How Google uses AI to optimize cooling in its data centers

- By applying DeepMind's machine learning to Google data centers, Google has managed to

- Reduced the amount of energy used for cooling by up to 40%.

- Thousands of sensors in the data center

- Every five minutes, Google's cloud-based AI pulls a snapshot of the data center cooling system

- Data is fed it into deep neural networks.

- The AI system then identifies which actions will minimize the energy consumption while satisfying a robust set of safety constraints.



Fun fact: Microsoft used the equivalent water volume of 2,560 Olympic size swimming to cool its data centers in 2022!

# AI in Your Business

## Benefits

- Cost Savings: AI can automate routine tasks, reducing labor costs.

- Improved Customer Experience: AI can provide personalized experiences and 24/7 customer service.

- Data Analysis: AI can analyze large amounts of data to provide insights.

- Increased Productivity: AI can take over mundane tasks, freeing up humans for more complex tasks.

# Deeplite Marketing Example

## Marketing campaign in less than half a day



**26 Page research paper in Arxiv**



**Website blog**



**LinkedIn post**

# Challenges of Implementing AI

- Training Data: Based on the application, sourcing training data can be difficult and expensive.

- Data Privacy: Businesses must ensure they respect customer data when using AI.

- Data security: Sharing your proprietary data with public AI could expose your secret sauce.

- Job Displacement: While AI can create jobs, it can also displace workers.

- Bias: Poorly designed AI may inadvertently lead to suspicions of discrimination.

- Implementation Costs: Developing and implementing AI can be expensive.

- Regulations: Governments globally are implementing "guardrails" to ensure AI is ethically developed and deployed

# What are employees saying?

## Study by SnapLogic

**81%** — Believe AI improves their overall performance at Work

**68%** — Want their employers to deploy more AI-based technology

**56%** — Stated they are already using AI to assist with their daily work

**89%** of employees believe that AI could support them in up to half of their workload

### Top benefits employees gain from using AI

- **61%** A more efficient and productive work day
- **49%** Improved decision making or faster time to insights
- **38%** Improved creativity
- **37%** Better collaboration between teams
- **35%** Better engagement with customers

**51%** — Believe AI helps them achieve a better Work/life balance

### What top three Workplace tasks do you think would benefit from AI?

- **43%** Understanding data and how trends and patterns can aid decision making
- **41%** Moving data from one place to another
- **41%** Accessing data residing in different places across the business

# Preparing Your Business for AI

- Start by identifying areas where AI could improve efficiency or customer service.

- Look to your vendors – what are they doing?

- Start small and get comfortable with it

- Invest in training and education to ensure your team understands AI.

- Develop a clear plan for implementing AI, including goals and timelines.

- Consider working with AI experts or consultants to ensure successful implementation.
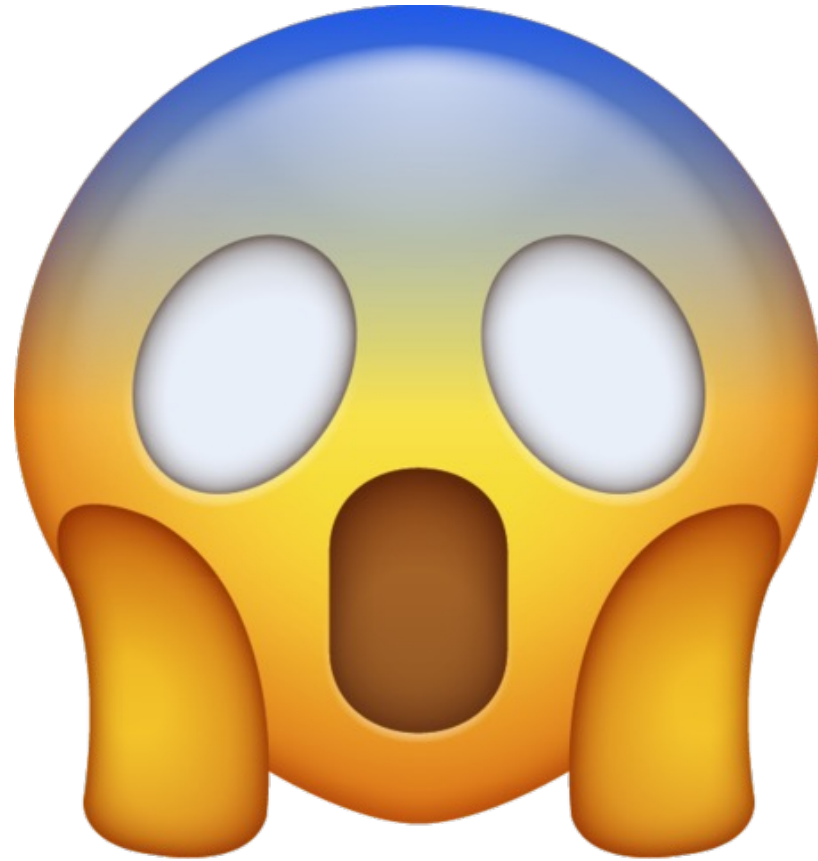
# Future of AI in Business

- AI will become increasingly integrated into business operations.

- It will become more ubiquitous

- Businesses that fail to adopt AI risk being left behind.

- AI will continue to create new job opportunities while displacing others.

- Ethical considerations will become increasingly important as AI becomes more prevalent.

# Ethical Considerations

- AI is a dual use technology

- Businesses have a responsibility to use AI ethically.

- This includes respecting customer data and privacy.

- Businesses must also consider the impact of AI on jobs and work to mitigate negative impacts.

- Transparency is key: Businesses should be clear about how and why they're using AI.

- Misuse will lead to over-regulation by governments

- AGI is closer than we think
- AI will have a profound impact on society
- Where there is data, there is an AI use case.  Data is gold.
- AI is becoming democratized – accessible to many
- Start small but make it a priority
- Look for areas of opportunity in your business
- Buy vs. build – look to your vendors
- Use AI responsibly and ethically

And let's hope the world doesn't end!

Confidential

# THANK YOU!

Please email me to let me know of a cool AI solution you've done!

# Sources

Source(s)

1. [Amazon's Secret to AI-Powered Product Recommendations](#)

2. [10 Ways Amazon Uses AI to Revolutionize E-Commerce in 2024](#)

3. [How Starbucks Leveraged AI Predictive Analytics for Personalized …](#)

4. [Leveraging AI for Personalized Customer Experiences. Lessons from …](#)

5. [How Starbucks is Revolutionizing Customer Relationship Management (CRM …](#)

6. [How Amex Helps You Protect Yourself Against Credit Card Fraud](#)

7. [American Express Prevents Fraud and Foils Cybercrime With NVIDIA AI …](#)

8. [Data centers are more energy efficient than ever - The Keyword](#)

9. [DeepMind AI reduces energy used for cooling Google data centers by 40%](#)

10. [Safety-first AI for autonomous data center cooling and industrial control](#)

11. [Artificial Intelligence at American Express – Two Current Use Cases](#)

12. [American Express Adopts NVIDIA AI to Help Prevent Fraud and Foil …](#)

13. [AI for humanity: How Starbucks plans to use technology to nurture the …](#)

14. [3 Ways Amazon Uses AI to Make Product Recommendations - Lineate](#)

15. [How Does Amazon Use Artificial Intelligence? Exploring AI-Powered …](#)

16. [Google's Use of AI to Manage Data Centers Enters a New Phase | Data …](#)