



New Frontiers

Richard Brooks
MAL2 – Spring 2025

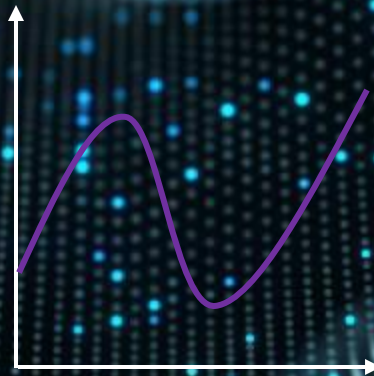
Credits: [MIT 6.S191 Spring 2025](#)

So far in MAL1+2

Data

- Signals
- Images
- Sensors

....



Functions Aproximators

Decision

- Prediction
- Predictions
- Actions

....

Power of Neural Networks

Universal Approximation Theorem

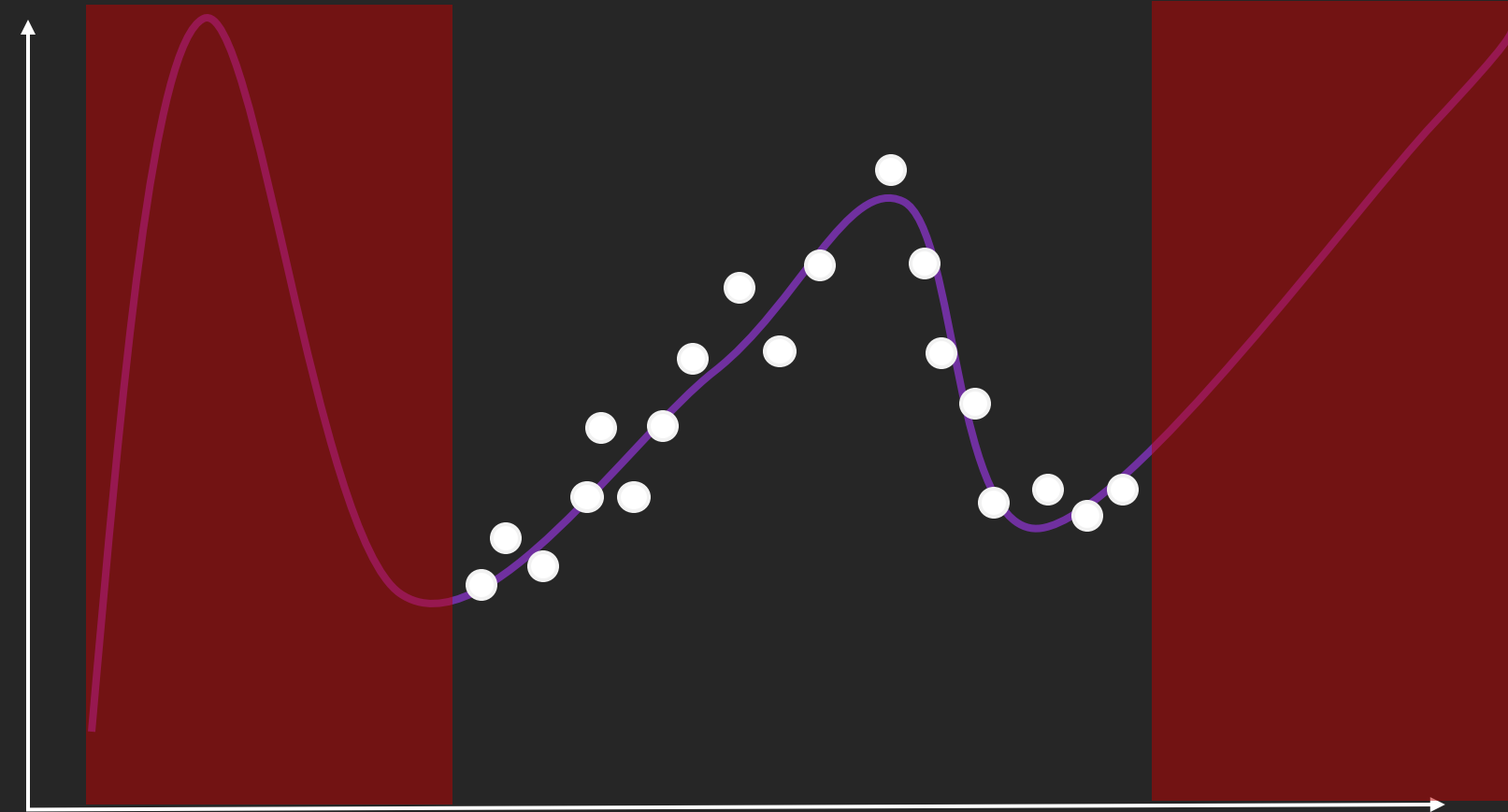
A feedforward network with a single layer is sufficient to approximate, to an arbitrary precision, any continuous function

Limitations



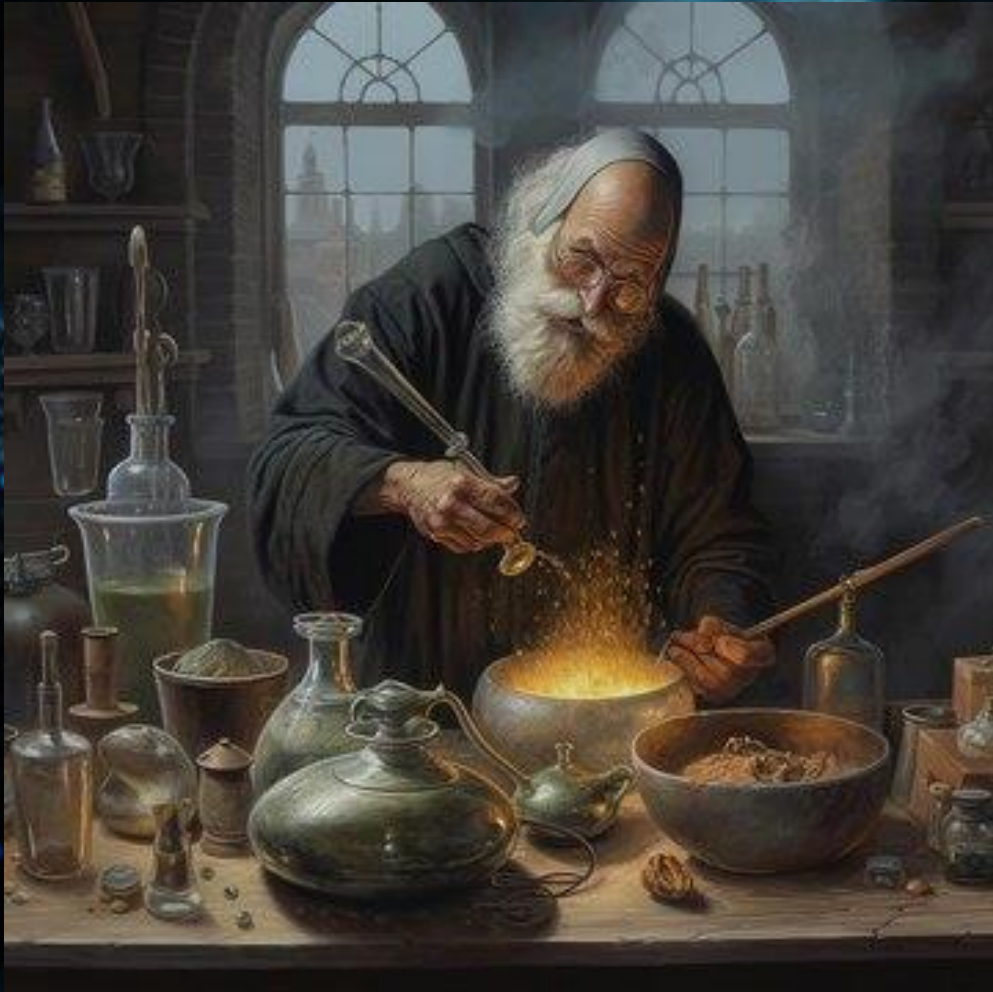
Neural networks are excellent function approximators
...when they have training data

Limitations



How do we know when our network doesn't know?

Deep Learning = Alchemy?



NN Failure

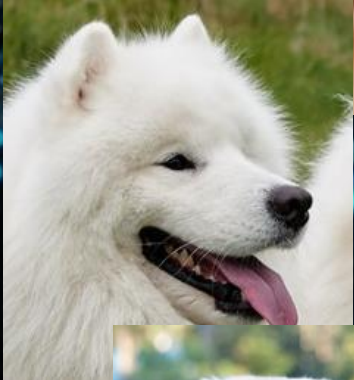


Train network to
colorize BW images.

Why could this be the case?

What Happens During Training

CNN



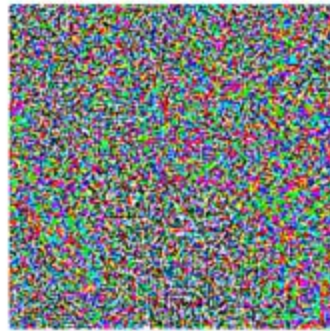
NN Failure



“panda”

57.7% confidence

+ .007 ×



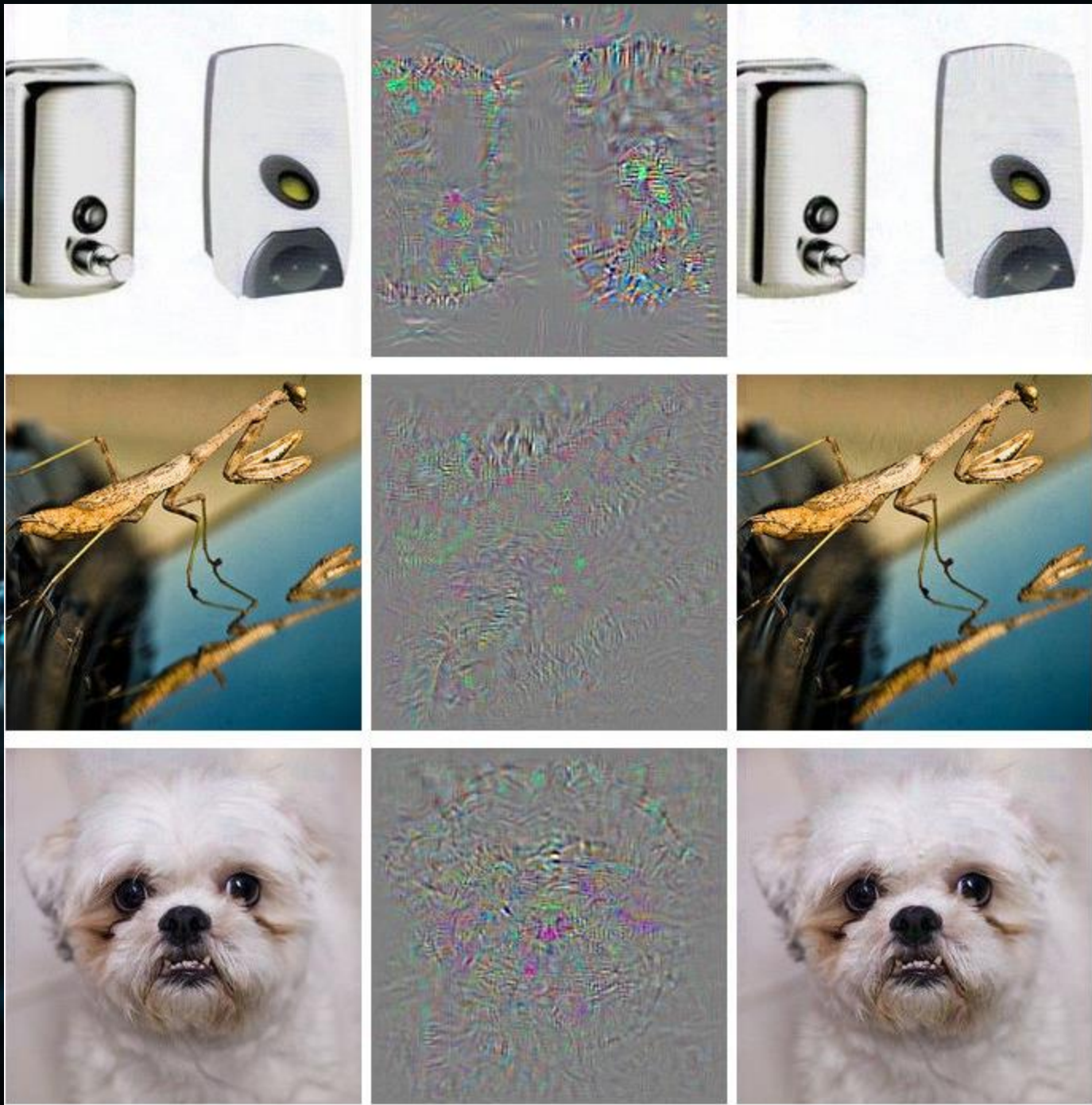
noise

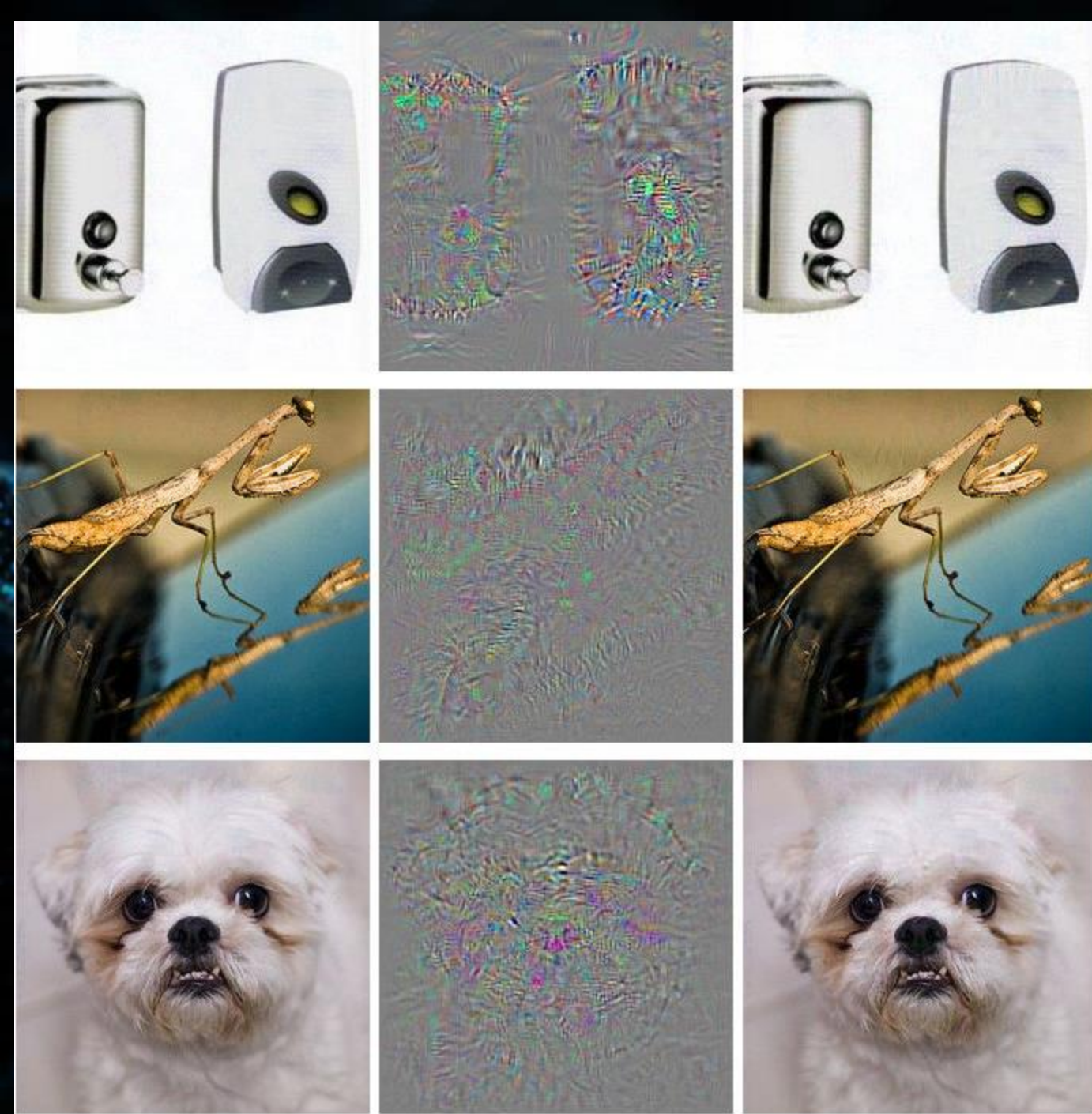
=



“gibbon”

99.3% confidence





These adversarial examples
were generated by minimizing
the following function with
respect to

$$\text{loss}(\hat{f}(\mathbf{x} + \mathbf{r}), l) + c \cdot |\mathbf{r}|$$

NN Limitations

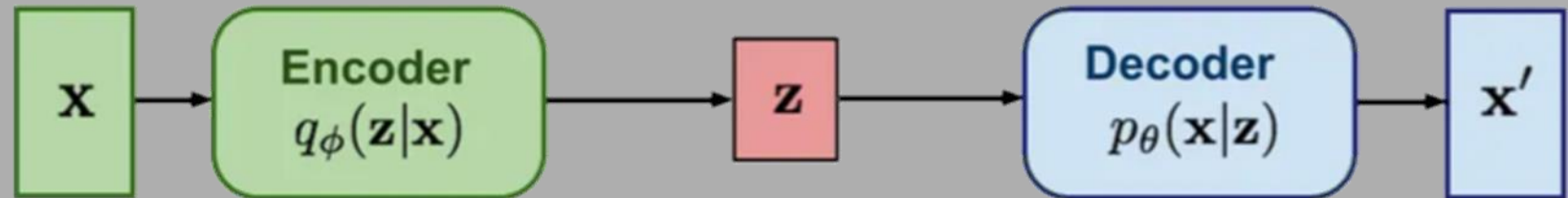
- Very **data hungry** (eg. often millions of examples)
- **Computationally intensive** to train and deploy (tractably requires GPUs)
- Easily fooled by **adversarial examples**
- Can be subject to **algorithmic bias**
- Poor at **representing uncertainty** (how do you know what the model knows?)
- Uninterpretable **black boxes**, difficult to trust
- Often require **expert knowledge** to design, fine tune architectures
- Difficult to **encode structure** and prior knowledge during learning
- **Extrapolation**: struggle to go beyond the data
- Hallucinations

Generative AI

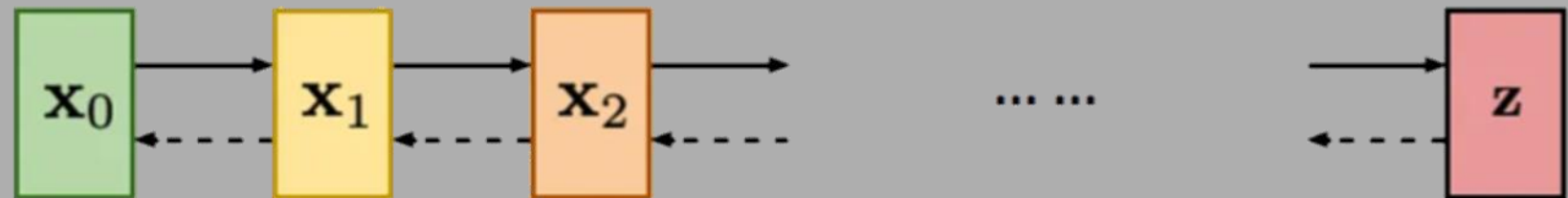
GAN: Adversarial training



VAE: maximize variational lower bound

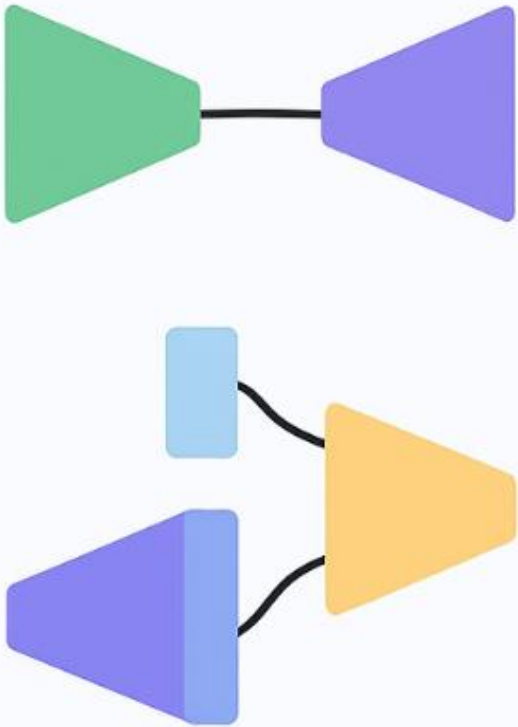


Diffusion models:
Gradually add Gaussian noise and then reverse

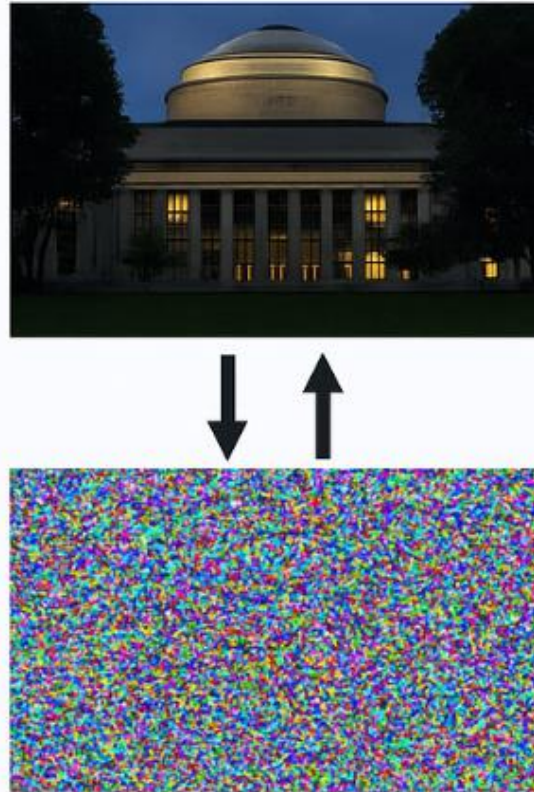


The Landscape of Generative Modeling

VAEs and GANs



Diffusion Models



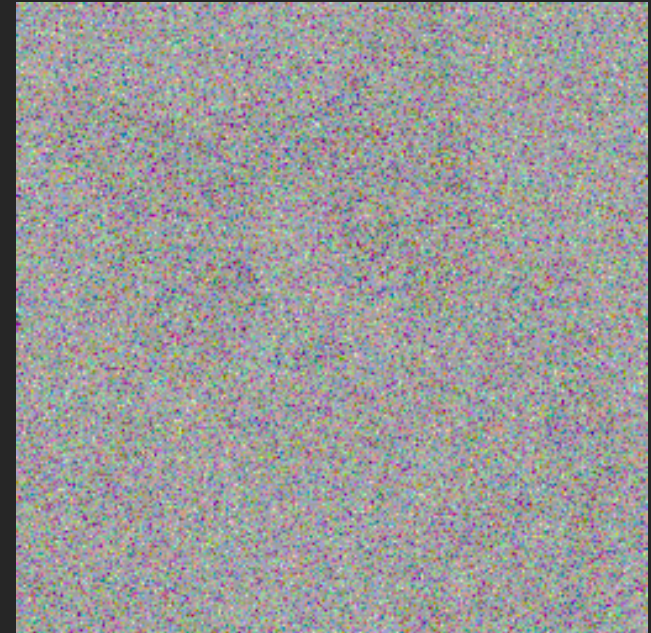
Text-to-Image



“Two cats doing research”

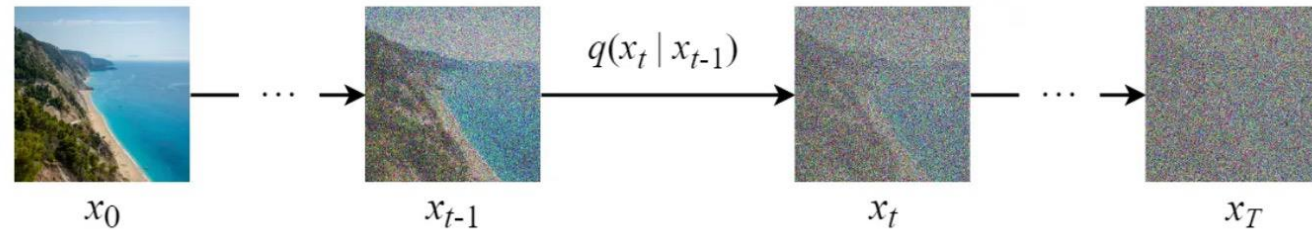
Diffusion Models

- Algorithm for training a diffusion model
 - **Forward process** — the process of adding noise to the image until it becomes a pure noise image.
 - **Reverse process** — the process where the model learns how to denoise the image gradually. This process starts with a complete noise image and should finish with a good looking image.
 - Reverse process makes up the **generative** model of the data.



Forward pass

- ① Forward Diffusion
 - ↳ add noise (Gaussian)
 - ↳ time steps
 - ↳ Markov Chain



Distribution of the
noised images

Output

Mean μ_t

Variance Σ_t

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I)$$

Notations:

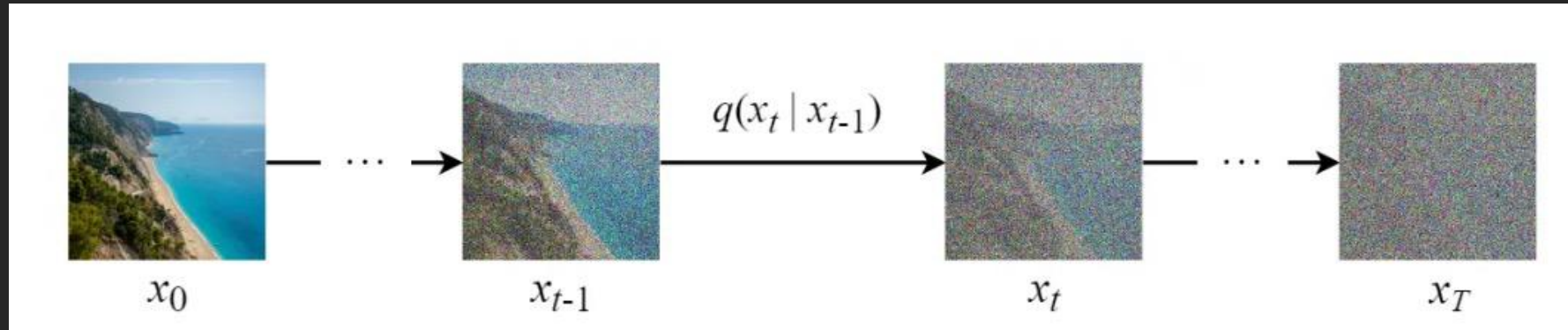
t : time step (from 0 to T)

x_0 : a data sampled from the real data distribution $q(x)$ (i.e. $x_0 \sim q(x)$)

β_t : variance schedule ($0 \leq \beta_t \leq 1$, and $\beta_0 = \text{small number}$, $\beta_T = \text{large number}$)

I : identity matrix

Reverse Denoising Process



so the main objective of the model here is to learn a function which predicts the mean and variance of each step to predict the reverse process.

② Reverse Diffusion
↳ Remove Noise
↳ CNN.V.Net
↳ MSE

Reverse Proces

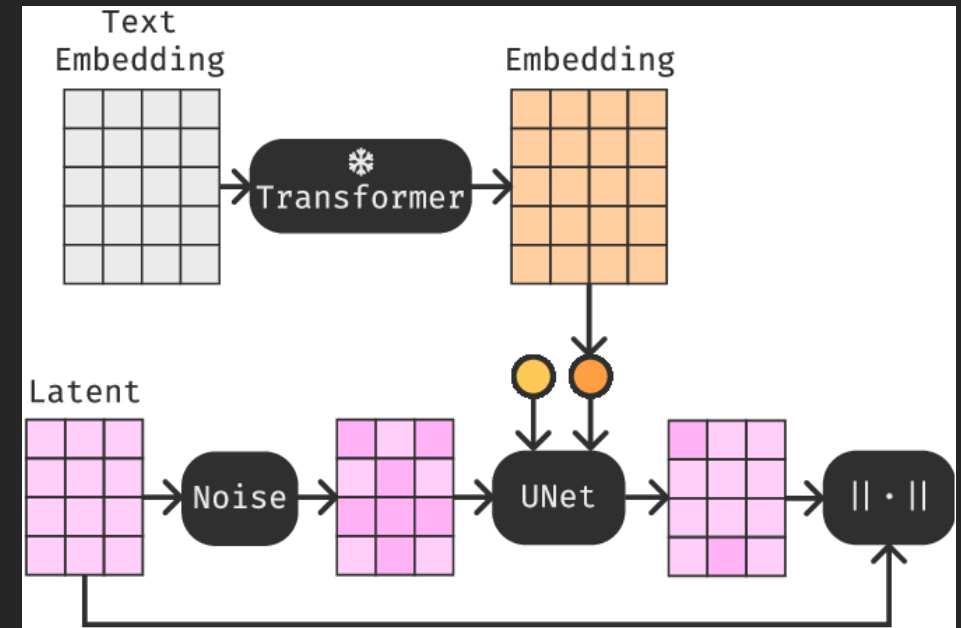
- The mean and variance in the reverse trajectory is learned, and the final function is used to undo the diffusion process.
- Final reverse process:

$$p(x_{t-1} \mid x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t))$$

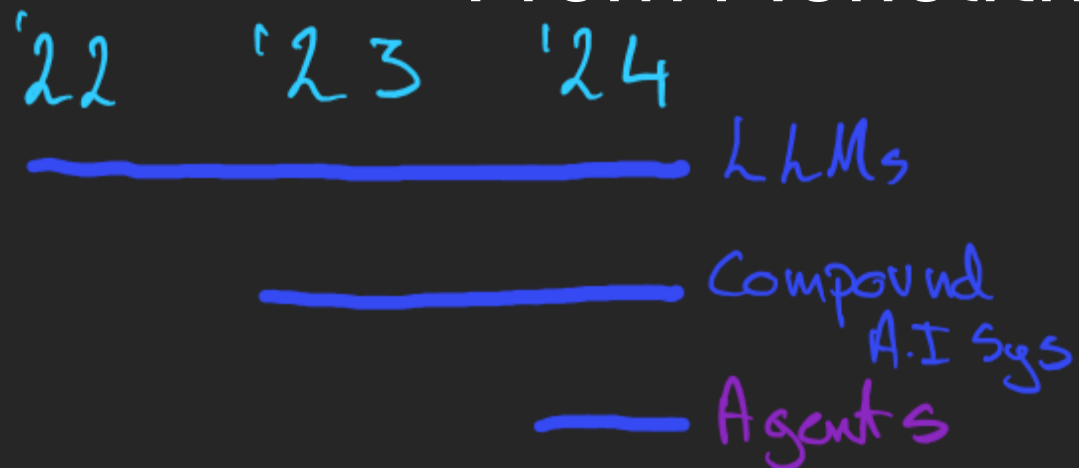
- Where:
- - $\mu(x_t, t)$ is the **mean**
- - $\Sigma(x_t, t)$ is the **variance**

Conditional Diffusion (Text-to-Image Guidance)

- Text Representation (Embeddings)
 - Text prompts need to be converted into numerical representations that capture their meaning.
- Conditioning the Denoising:
 - The core idea is that the text embedding is provided as additional input to the U-Net during the reverse diffusion (denoising) process at each time step t .
- Guided Noise Prediction:
 - The U-Net now learns to predict the noise not only based on the noisy image x_t but also conditioned on the provided text embedding. It learns to remove noise in a way that steers the image formation towards something that matches the text description.



From Monolithic Models to AI Agents

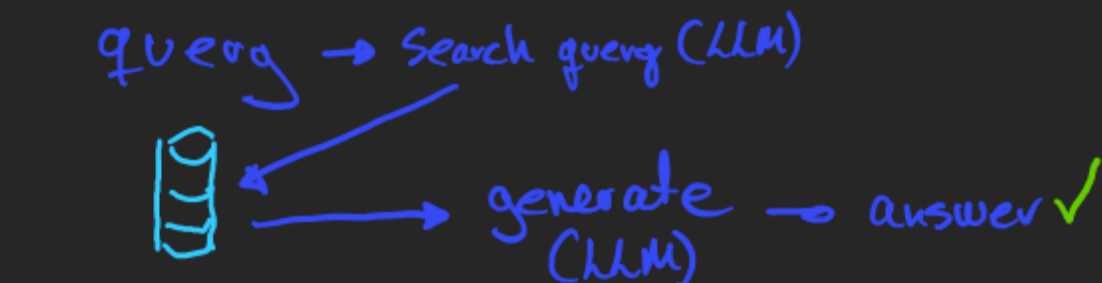


From models \rightarrow Compound sys. ^{RAG}

- ⊕ limited knowledge
- ⊕ hard to adapt

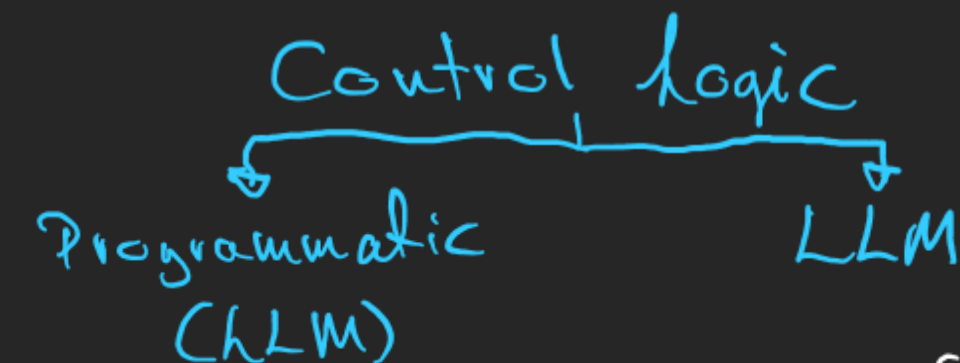
☀ Vacation days:

Query \rightarrow generate (LLM) \rightarrow answer X



System Design

- ⊕ Modular
- ⊕ Easier to adopt

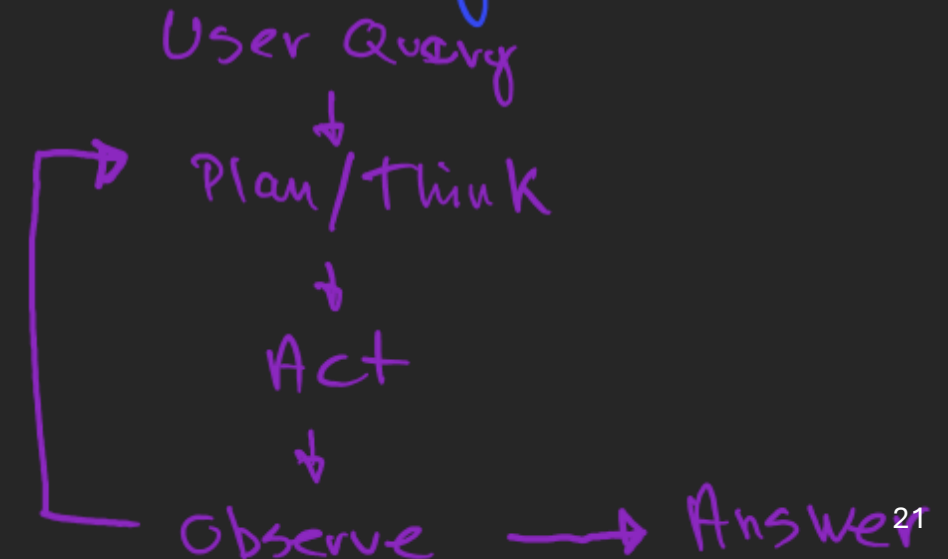
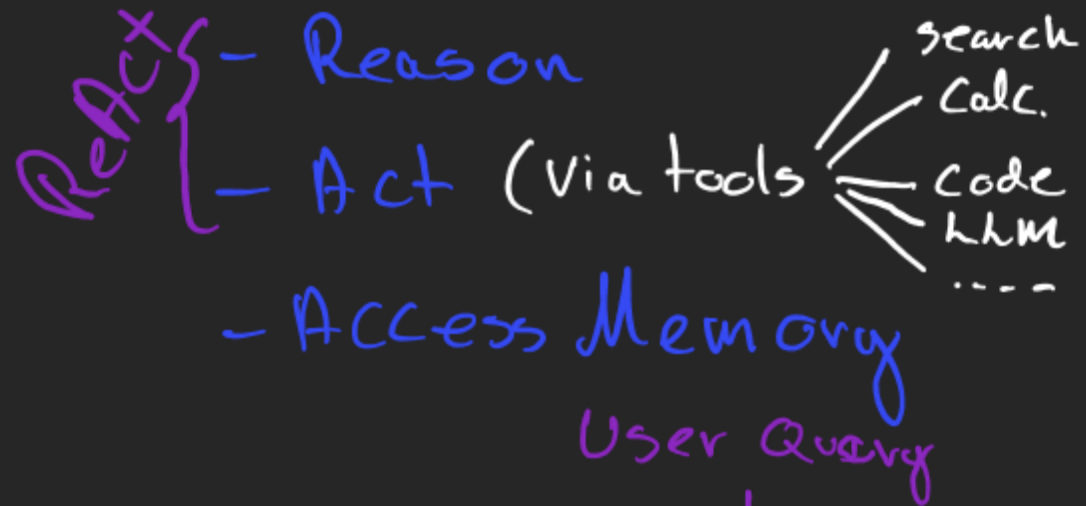
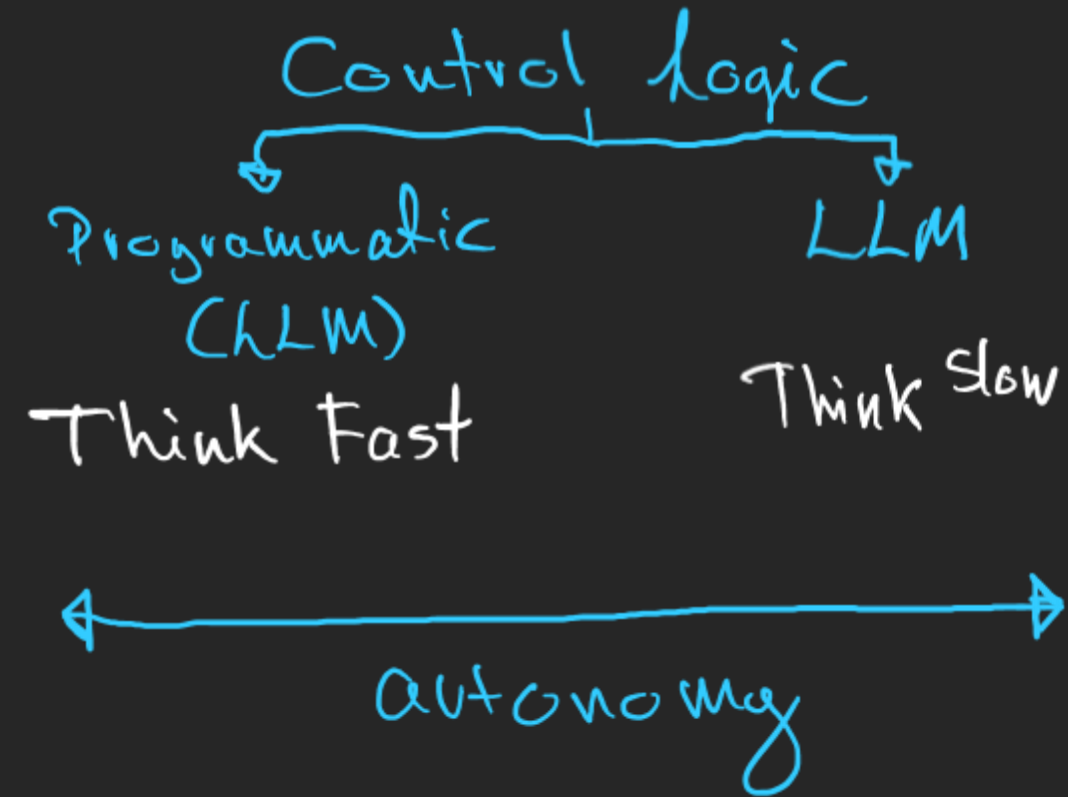


Think Fast

Think Slow

LLM Agents

LLM Agents



Non=Agentic Workflow (Zero-Shot)

One-Pass Execution
- Completes a task in one go without iteration.

No Adjustments
- No revisions or refinements.

Agentic Workflow

Step-by-Step Execution
- Breaks tasks into stages (e.g., planning, research, drafting, revising).

Feedback Loop
- AI iterates based on human guidance.

Truly Autonomous AI Agent

Fully Independent AI
- AI determines steps, tools, and iterates without human involvement

Adaptive Decision-Making
- AI adjusts workflow dynamically to achieve the best outcome.

Design Patterns

1. Reflections

- AI reviews and improves its own output by analysing and refining responses

2. Tool Use

- Agents leverage external tools like web search or code execution

3. Planning and Reasoning

- AI determines optimal steps and selects tools to execute tasks efficiently

4. Multi-Agent Systems

- Multiple AI model with specialized roles
 1. Sequential
 2. Hierarchal Pattern
 3. Parallel Systems
 4. Asynchronous Systems

The compas algorithm

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

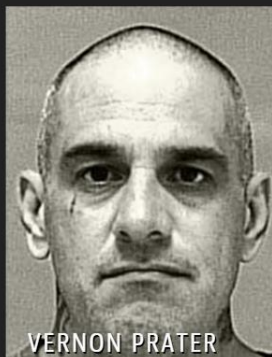
May 23, 2016



10. Fairness and Non-Discrimination

AI actors should promote social justice, fairness, and non-discrimination while taking an inclusive approach to ensure AI's benefits are accessible to all.

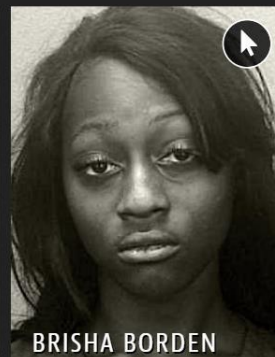
Two Petty Theft Arrests



VERNON PRATER

LOW RISK

3



BRISHA BORDEN

HIGH RISK

8

Prior offences

4 juvenile misdemeanors

Subsequent offences

None

6. Transparency and Explainability

The ethical deployment of AI systems depends on their transparency & explainability (T&E). The level of T&E should be appropriate to the context, as there may be tensions between T&E and other principles such as privacy, safety and security.

7. Human Oversight and Determination

Member States should ensure that AI systems do not displace ultimate human responsibility and accountability.

Prior offences

2 armed robberies

1 attempted armed robbery

Subsequent offences

1 grand theft

Uber Self-driving Car Fatality

Uber self-driving car kills pedestrian in first fatal autonomous crash

Driver Charged in Uber's Fatal 2018 Autonomous Car Crash

Uber self-driving car test driver pleads guilty to endangerment in pedestrian death case

Driverless cars are mostly safer than humans – but worse at turns

2. Safety and Security

Unwanted harms (safety risks) as well as vulnerabilities to attack (security risks) should be avoided and addressed by AI actors.

5. Responsibility and Accountability

AI systems should be auditable and traceable. There should be oversight, impact assessment, audit and due diligence mechanisms in place to avoid conflicts with human rights norms and threats to environmental wellbeing.

Clearview Ai Facial Recognition

The Secretive Company That Might End Privacy as We Know It

A little-known start-up helps law enforcement match photos of unknown people to their online images — and “might lead to a dystopian future or something,” a backer says.

Clearview AI— Controversial Facial Recognition Firm— Fined \$33 Million For 'Illegal Database'

AI can now Identify People as Gay or Straight from their Photo

3. Right to Privacy and Data Protection

Privacy must be protected and promoted throughout the AI lifecycle. Adequate data protection frameworks should also be established.

1. Proportionality and Do No Harm

The use of AI systems must not go beyond what is necessary to achieve a legitimate aim. Risk assessment should be used to prevent harms which may result from such uses.

Chatbots

5. Responsibility and Accountability

AI systems should be auditable and traceable. There should be oversight, impact assessment, audit and due diligence mechanisms in place to avoid conflicts with human rights norms and threats to environmental wellbeing.

Can A.I. Be Blamed for a Teen's Suicide?

Snapchat-bot giver børn detaljerede råd om at selvskade

Man ends his life after an AI chatbot 'encouraged' him to sacrifice himself to stop climate change

1. Proportionality and Do No Harm

The use of AI systems must not go beyond what is necessary to achieve a legitimate aim. Risk assessment should be used to prevent harms which may result from such uses.

Sustainability

Artificial intelligence technology behind ChatGPT was built in Iowa — with a lot of water

ChatGPT And Generative AI Innovations Are Creating Sustainability Havoc

A.I. Could Soon Need as Much Electricity as an Entire Country

8. Sustainability

AI technologies should be assessed against their impacts on 'sustainability', understood as a set of constantly evolving goals including those set out in the UN's Sustainable Development Goals.

A photograph of two people in business attire shaking hands, symbolizing agreement or partnership. The image is overlaid with a semi-transparent dark blue rectangle containing text.

1. Proportionality and Do No Harm

The use of AI systems must not go beyond what is necessary to achieve a legitimate aim. Risk assessment should be used to prevent harms which may result from such uses.

A close-up of a person wearing a VR headset, looking intently at the virtual environment. The image is overlaid with a semi-transparent dark blue rectangle containing text.

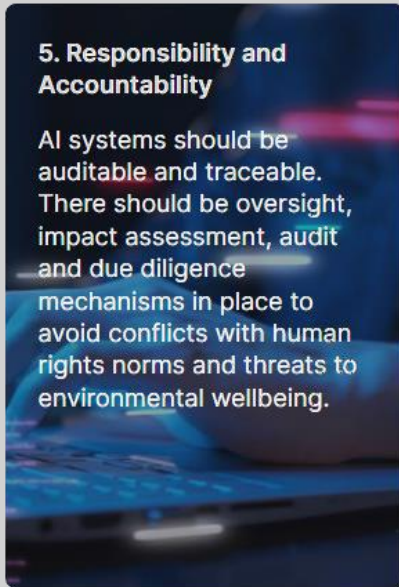
3. Right to Privacy and Data Protection

Privacy must be protected and promoted throughout the AI lifecycle. Adequate data protection frameworks should also be established.

A person looking at a tablet with various digital overlays and icons, representing multi-stakeholder governance. The image is overlaid with a semi-transparent dark blue rectangle containing text.

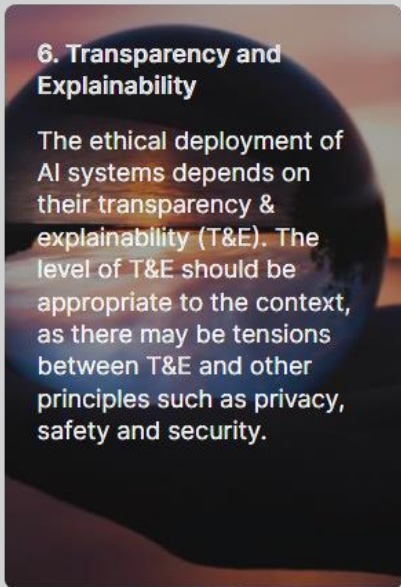
4. Multi-stakeholder and Adaptive Governance & Collaboration

International law & national sovereignty must be respected in the use of data. Additionally, participation of diverse stakeholders is necessary for inclusive approaches to AI governance.

A person looking at a tablet with various digital overlays and icons, representing responsibility and accountability. The image is overlaid with a semi-transparent dark blue rectangle containing text.

5. Responsibility and Accountability

AI systems should be auditable and traceable. There should be oversight, impact assessment, audit and due diligence mechanisms in place to avoid conflicts with human rights norms and threats to environmental wellbeing.

A person looking at a tablet with various digital overlays and icons, representing transparency and explainability. The image is overlaid with a semi-transparent dark blue rectangle containing text.

6. Transparency and Explainability

The ethical deployment of AI systems depends on their transparency & explainability (T&E). The level of T&E should be appropriate to the context, as there may be tensions between T&E and other principles such as privacy, safety and security.

A person looking at a tablet with various digital overlays and icons, representing safety and security. The image is overlaid with a semi-transparent dark blue rectangle containing text.

2. Safety and Security

Unwanted harms (safety risks) as well as vulnerabilities to attack (security risks) should be avoided and addressed by AI actors.

A person looking at a tablet with various digital overlays and icons, representing human oversight and determination. The image is overlaid with a semi-transparent dark blue rectangle containing text.

7. Human Oversight and Determination

Member States should ensure that AI systems do not displace ultimate human responsibility and accountability.

A person looking at a tablet with various digital overlays and icons, representing sustainability. The image is overlaid with a semi-transparent dark blue rectangle containing text.

8. Sustainability

AI technologies should be assessed against their impacts on 'sustainability', understood as a set of constantly evolving goals including those set out in the UN's Sustainable Development Goals.

A person looking at a tablet with various digital overlays and icons, representing awareness and literacy. The image is overlaid with a semi-transparent dark blue rectangle containing text.

9. Awareness & Literacy

Public understanding of AI and data should be promoted through open & accessible education, civic engagement, digital skills & AI ethics training, media & information literacy.

A colorful illustration featuring various people, including one in a wheelchair, surrounded by AI-related symbols like gears, a brain, and a smartphone. The image is overlaid with a semi-transparent dark blue rectangle containing text.

10. Fairness and Non-Discrimination

AI actors should promote social justice, fairness, and non-discrimination while taking an inclusive approach to ensure AI's benefits are accessible to all.

Make an inspirational quote to end my advanced machine learning course MAL2 with.



Every model we build expands not just our knowledge, but our ability to see patterns and possibilities in the world around us. As you move forward, remember: the true power of machine learning isn't in algorithms – it's in the questions we dare to ask and the insights we strive to reveal.

Keep learning, stay curious, and continue transforming data into discovery.

Best of luck!