

# 1. Symbolic / Rule-based AI ("Expert Systems")

- GOFAI (Good Old-Fashioned AI)

- Uses explicit symbols, rules, logic to perform reasoning

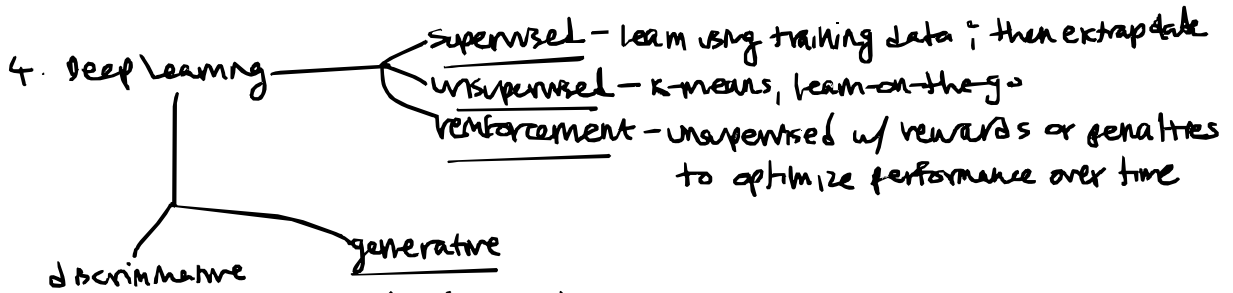
# 2. Modern AI are built off neural networks.

- Model AI off human brain.

- Nodes in neural network are associated w/ dynamically adaptable weights, which change w/ inputs to system

# 3. Shallow learning vs. Deep Learning

few layers      many layers  
in neural networks



- Model  $P(Y|X)$  -  
prob. of label for input  
- logistic regression,  
decision trees  
- Use for classification

- Use for creation  
- Model  $P(Y, X)$  -  
joint prob. that data &  
labels occur together  
- GANs, Naive Bayes

Ask yourself "Is model creating data?" If yes, generative.  
Otherwise, discriminative.

# 5. Features (input/activation layer)

→ hidden layers → output layer

w/ deep learning, we have hundreds/thousands of hidden layers, w/ dynamically adjusting weights.

6. Word embeddings - Representation of text where words or phrases from vocab. are mapped to vectors of real #s

└ words w/ similar meanings are closer together.  
Those w/ diff meanings are further apart.

7. Transformer models - processes word embeddings

└ use attention mechanisms to understand rels b/w words in a sentence

8. Gen. AI systems are actually predictive systems.

└ "stochastic parrots" - They repeat things w/out true comprehension.

9. Limits of Algorithmic systems

- Related to 3 key attributes

- └ data-driven
- └ based on statistical models
- └ use dynamic w/ development

- Data-Driven

\* Data is not reality. It is a representation / model of reality to hone in on the aspects of reality closest to the problem we're solving

\* Data requires quantification.

\* Any representation of data is a choice - and can leave things out.

- Stat. Models

\* Fundamentally correlational

\* Pattern recognition engines

\* No understanding of causal relationships

10. AI Accuracy

- Capability: Meets AI's defined objective fcn.

└ if a classifier is producing wrong classifications

- Alignment: Meets intended human objective

└ includes things like scope for LLMs — Don't use customer service bot for medical diagnosis.  
└ toxicity

- Robustness: May be accurate at launch --- but does it stay accurate?

11. Overfitting - Model has basically memorized the training data instead of learning underlying patterns.  
So, for anything outside training data, it will perform horribly.  
TLDR: overfitting ruins generalization.

## 12. Distinctive Generative AI Accuracy Issues

- Unpredictable applications
- Hallucinations
  - ↳ Not tied to existing examples in dataset
- Concentration risk of found. models
  - ↳ You end up getting just a few companies owning & running actual found. models, just like w/ hyperscalers.
- Need to have room for error? randomness b/c generative AI has to be creative - b/c it's creating data
- PEOKAC (Problem exists w/o keyboard and chair)
  - ↳ People = source of error too but still an integral part of sys!

## 13. Transparency — Awareness: Is an AI system being used?

- ↳ Disclosure: What are the features/attributes?  
What are the model weights, tuning params, guardrails, etc?
- ↳ Interpretation or Explanation  
Why did the system produce the result that it did?

## 14. Improve transparency thru

- algorithmic audits — 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> party
- model/system cards — vendor publishes notes on model + how it works + how built

## 15. Interpretable - linear models subject to direct understanding

- ↳ Deep learning / LLMs = not interpretable b/c of many hidden layers

## 16. Explainable - finding interpretable models consistent w/ the output

## 17. Limits of Interpretability / Explainability

- Technical Illiteracy
- Not a description of the actual "black box"
  - ↳ just an approximation

18. You can ask a LLM to explain itself... but reliability is not guaranteed.

- Buyer beware - LLM fundamentally non-deterministic, so it can be and change response to prompt for explanation every time

19. Common Means of Explainability — using technical tools to describe model's features & workings

Compare to simpler model w/ similar objective fun. — Companion model to predict & explain complex main model's outputs

20. Label Choice Bias — Gap b/w ideal (what we want) vs actual (what it actually does)

- ↳ e.g., AI models diagnosing illnesses or reducing cost of care by disregarding most vulnerable

21. Bias

### Historical Bias

People's biased views and decisions infect the data being fed into the model.

- ↳ e.g., Amazon resume screener favored men b/c patriarchy has historically favored men over women in professional settings — which

### Evaluation Bias

Mechanisms used to evaluate AI themselves have & introduce bias.

- ↳ All facial recognition systems trained off / optimizing from common CV database — which was massively biased

### Representation Bias

Training data doesn't reflect reality; underrep. of certain groups

- ↳ minorities — Facial recog. systems less accurate for black female faces
- ↳ poor

### Aggregation Bias — Don't be fully blind to differences.

Data from diff groups are combined w/out accounting for underlying patterns or relationships

City A: 1 million people, \$50 avg.

City B: 100 people, \$40 avg

AB → If you say avg. income =  $\frac{50 \times 30}{2} = 40$

0/c actual avg =  $\frac{50(1e6) + 30(100)}{1e6 + 100}$

Avg has to take # of people in each city into account.

22. Bias doesn't just come from data - also from human choices in creating & deploying the model.
23. Deepfakes - identity theft, where bad actor uses AI system to pretend to be someone else in order to fool/exploit people
24. Synthetic data - using AI to create data that is then used to train other AI
25. In the US, generative AI outputs are not protected by copyright.
- Copyright only protects human authorship.
  - Beijing court disagrees - beginning of a different tack in China
  - Tension b/w protecting copyright/IP vs. promoting AI innovation
26. Licensed/synthetic/public-domain data → Don't have to worry about copyright restrictions
27. Privacy Lifecycle
- Collection/Creation — what data to gather + how  
— Prioritize data minimization (only gather what's needed)
  - Aggregation & Analysis — data combined + processed
  - Use — data used to power apps/services
  - Storage — Is it safely stored? Encryption at rest/backups/access control
  - Distribution — Who is it shared with?
28. Metaprompting - "Here are the rules & patterns to follow when answering questions like this"
- "You are ...." - giving a persona as guideline for how to behave