



---

© SupportVectors. All rights reserved.

This document is the intellectual property of SupportVectors, and part of its training material. Only the participants in SupportVectors workshops are allowed to study the document for educational purposes currently, but is prohibited from copying or using it for any other purposes without written permission. These documents are chapters and sections from Asif Qamar's textbook that he is writing on Data Science. So we request you to not circulate the material to others.

---

## Summary of Week 8: AI Agents

### Table of Content

- **Fine-Tuning Models in the Context of AI Agents**
- **What is Model Training ?**
  - An Example of Ice Cream Sales
  - An Example of Cows and Ducks
- **The Universal Approximator Theorem**
- **The Neural Network Training Loop**
  - Model Mechanics and Non-Linearity
  - The Four Steps of Model Training
  - Null Hypothesis
- **Transfer Learning**
  - The Challenge of Scale and The Neural Scaling Laws
  - The Solution: Foundational Models and Foundational Ability
- **Training Foundational Abilities through Surrogate Tasks**
  - Surrogate Task: The Forcing Function
    - The Example of BERT (Guess the Masked Word)
  - Inducing Foundational Ability
    - The Mandarin Child Analogy
  - The Machine's Learning

- Transfer to Last-Mile Problems
- **Quantifying Loss and The Enterprise Use Case**
  - Loss Function for the Surrogate Task (BERT)
  - Transfer Learning in an Enterprise Context
  - Conclusion on Pricing Models
- **Modern Model Training: The Two-Step Journey of Transfer Learning**
  - The Unifying Concept: Transfer Learning
- **The Necessity of Fine-Tuning AI Agents**
- **Why Fine-Tuning Beats Prompt Engineering ?**
  - In-Context Learning vs. Deep Habit
  - Specialized Responsibility
  - The Last Mile Problem
- **Frontier Models vs. Fine-Tuning: A Cost-Benefit Analysis**
  - The Fine-Tuning Spectrum: Sizing the Brain
  - The Analogy of Specialization
  - The Escalation Ladder for Application Development
- **RAG vs. Fine-Tuning: A False Dichotomy**
- **Distinguishing Fine-Tuning from Knowledge Augmentation**
  - Why Fine-Tuning is for Behavior, Not Knowledge
  - The Fine-Tuning Stage in the Escalation Ladder
- **LLMs vs. Fine-Tuned Models**
  - Arguments for Large Frontier Models (LLMs)
  - Arguments for Fine-Tuned Models (SLMs / Web of Models)
  - The Crucial Trade-Offs
- **A2A Conceptual Deep Dive**
  - The Distinction: Client vs. Agent
- **A2A vs. Single-Process Frameworks (CrewAI, LangGraph)**
- **Interoperability: A2A, ADK, and MCP**
  - Agent Discovery: Peer-to-Peer vs. Centralized
- **The Future of A2A: Standardization and Security**
- **Advanced Agentic Concepts**
  - New Framework: Microsoft's Agent-FL for Fine-Tuning
  - The Foundation of Tool Calling
  - Agent Design Pattern: Planning and Execution
  - Agent Design Pattern: Re-planning
  - Agent Design Pattern: Reflection
  - Example: The "Smart Home" Reasoning Project

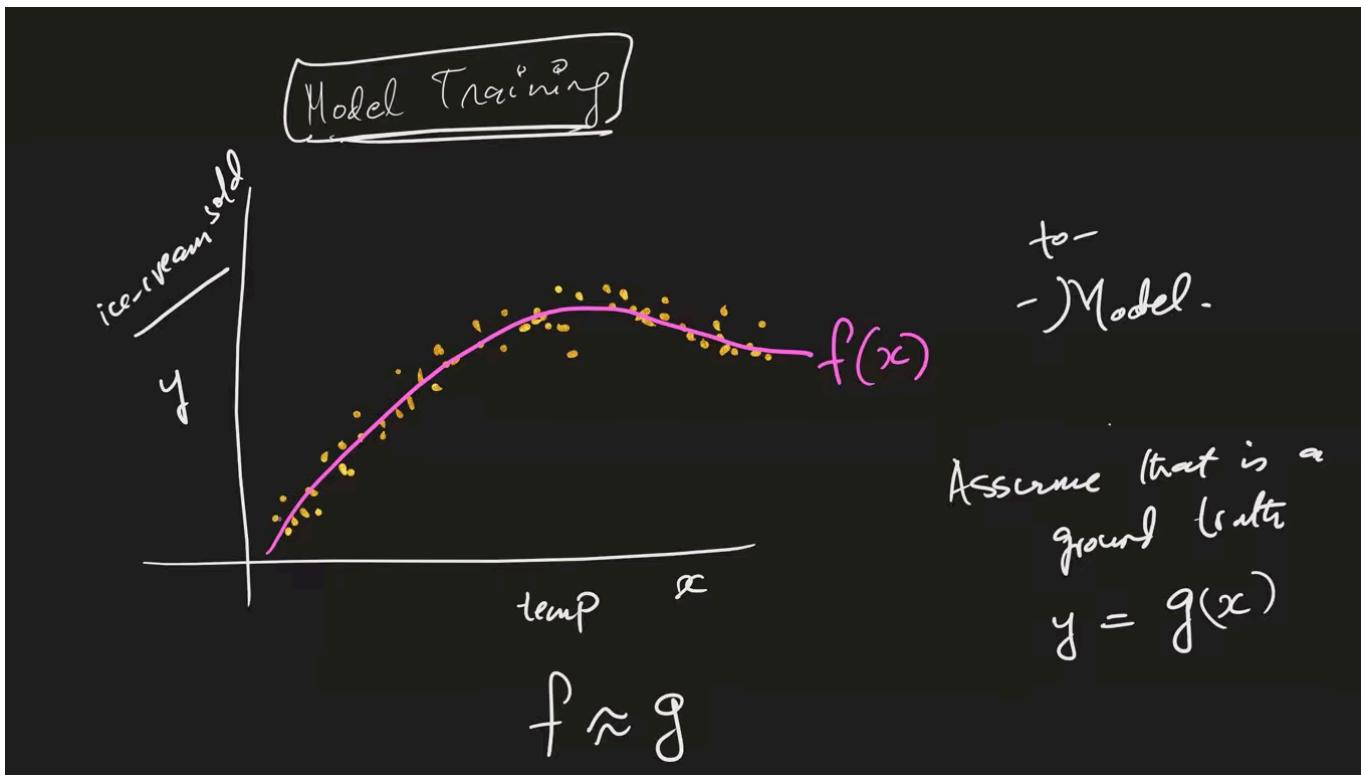
- **Practical Career and Skill Discussion**
  - Real-World Adoption of Complex Patterns
  - Key Skill Sets for the Evolving AI Field
- **The Case Against Fine-Tuning for Knowledge**
- **Knowledge Obsolescence and Cost**
- **Risks: Catastrophic Forgetting and Mode Collapse**
- **Talent and Practicality (The Escalation Ladder)**
- **Thermodynamic Disconnect: Knowledge Entropy vs. Model Capacity**
- **Parameter Limits and Grokking (Lossy Compression)**
- **Data Provenance and Legal Requirements**
- **The Fundamental Dichotomy and System Architecture**
  - The Fundamental Dichotomy:
  - What is Fine-Tuning?
  - The Maturity Escalation Ladder
- **Debate: Frontier Models vs. Fine-Tuned Local Models**
- **Project: The Patient Advocacy Agent**
  - Motivation and Narrative
  - Core Rule: Do Not Play Doctor
  - Medical Framework: SOAP
  - Required Outputs
  - Interface and Technical Requirements
  - Scope and Common Sense
- **Key Insights**

## Fine-Tuning Models in the Context of AI Agents

The session introduces **fine-tuning of models**, specifically within the context of **AI Agents**. It is emphasized that fine-tuning is a practical skill, likened to playing tennis, that requires hands-on practice for proficiency. To understand fine-tuning, one must first grasp the core concept of **model training** and the purpose it serves in machine learning.

---

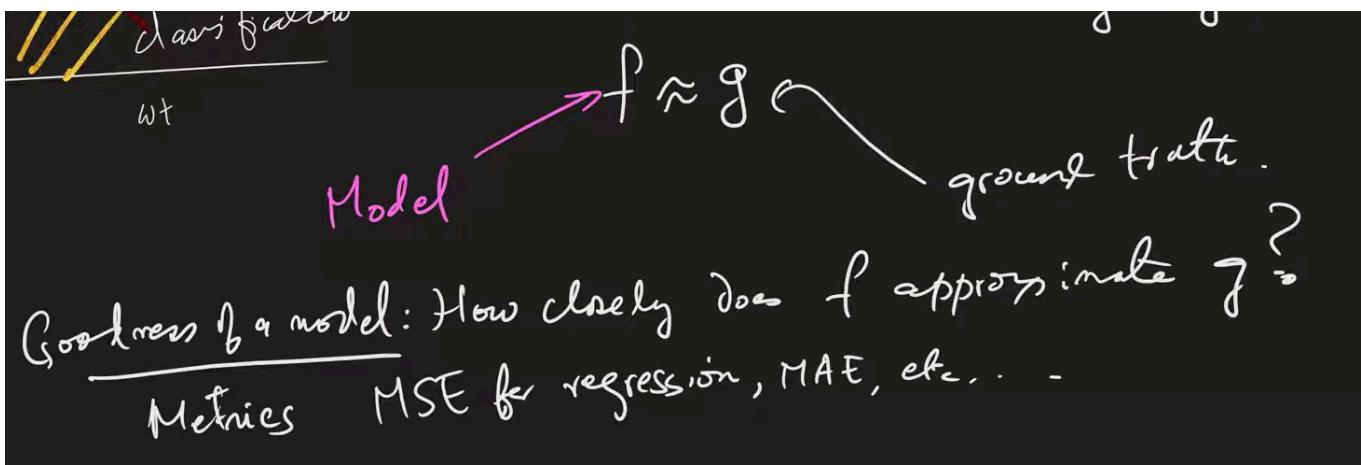
## What is Model Training ?



Model training is the process of **modeling** reality, which means searching for and finding an effective function,  $f(x)$ . It is assumed that observed empirical data is produced by an unknown, underlying **ground truth function**,  $y = g(x)$ . This function describes the true relationship between variables. Since the ground truth  $g(x)$  is unknown, the goal of all machine learning is to create a model,  $f(x)$ , that is a **close approximation** to  $g(x)$ , or  $f \approx g$ . The core idea is captured by the phrase: "**All models are wrong, but some models are useful.**" Models are technically "wrong" because they are merely approximations and not the true, unknown ground truth. However, they are useful if they are a **good approximation** of the underlying reality.

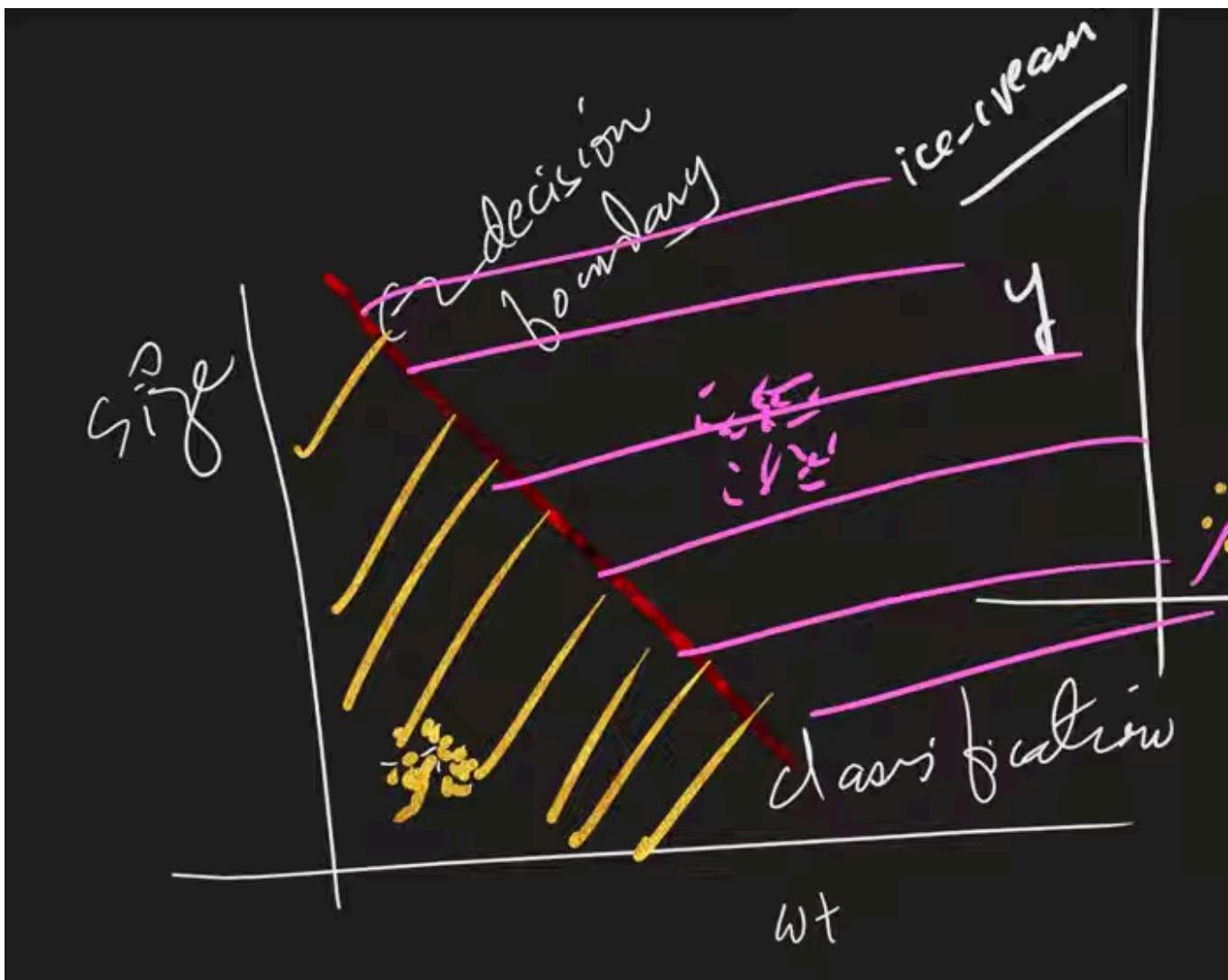
## An Example of Ice Cream Sales

An entrepreneur wants to predict daily ice cream sales (Ice Cream Sold,  $y$ -axis) based on the temperature ( $x$ -axis) to avoid waste or missed opportunity. We assume that the data is produced by  $y = g(x)$ , which is a ground truth function. The model must find a function  $f(x)$  that accurately predicts sales from temperature, allowing the entrepreneur to buy the right stock from the wholesaler. In other words, we need to find a function  $f(x)$ , that is a **close approximation** to  $g(x)$ , or  $f \approx g$ . The model suffers from **unknown unknowns** (e.g., wind, rain, holidays) that affect sales but aren't included as variables in the data, making the model an approximation.



For regression, the approximation's quality is measured using metrics like **Mean Squared Error (MSE)** and **Mean Absolute Error (MAE)**.

## An Example of Cows and Ducks



Consider a case where the goal is to classify animals based on features like weight and size. Ducks are small and feathery, while cows are big and heavy. To model this, the machine learning algorithm draws a **decision boundary** that separates the data points (cows from

ducks). For classification, the error rate is measured using metrics like **accuracy**, **precision**, **recall**, and **F1 score**.

The principles governing these basic machine learning problems apply universally to more complex systems, including Large Language Models. Even fundamental physical theories, like **Newton's Law of Gravity** (the observation of apples falling), are effective, useful, but ultimately approximate models of a deeper, unknown truth.

To train a model is to find a function  
 $f \approx g$ .  
that is an effective  
approximator of  $g$ .

\* No free lunch theorem

To train a model is therefore to find a function  $f$  that is an **effective approximator of  $g$** .

### Note

## The No Free Lunch Theorem

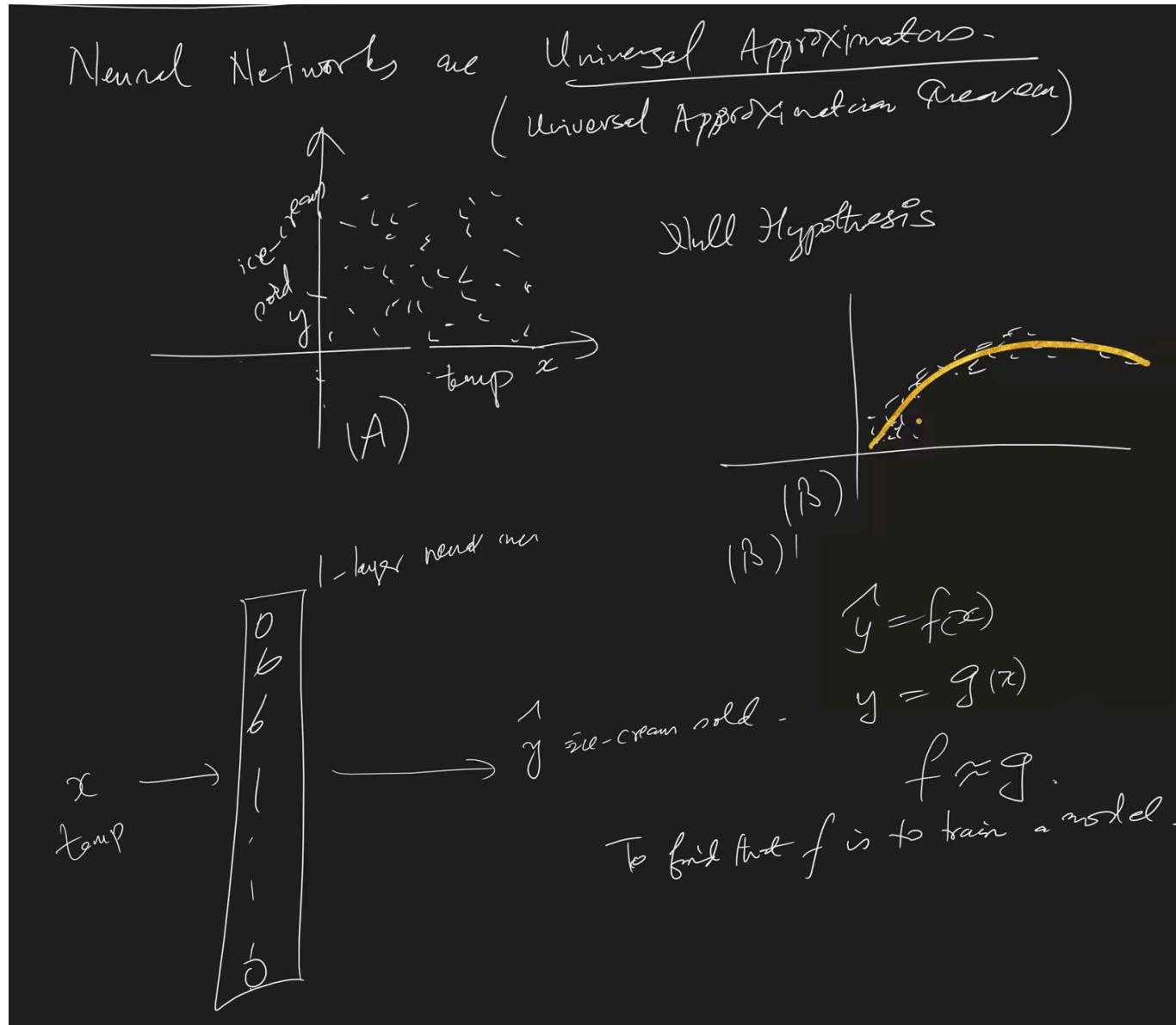
The theorem states that **no one particular kind of algorithm is necessarily better than all other algorithms** across all possible use cases and domains. The effectiveness of an algorithm is specific to its application. While one algorithm may excel in a particular domain, it is not "forever effective" or always the best choice in every situation.

## The Universal Approximator Theorem

It states that a standard feedforward neural network with a **single hidden layer** and a sufficient number of neurons can effectively approximate **any continuous function** to an arbitrary degree of accuracy. This is crucial because **real-world relationship** is expressed quantitatively through functions. When we observe phenomena, like an apple falling (governed by a function relating gravity and mass) or ice cream sales (governed by a relationship with temperature) we are seeing the results of underlying, often **non-linear**, functions. Since neural networks are

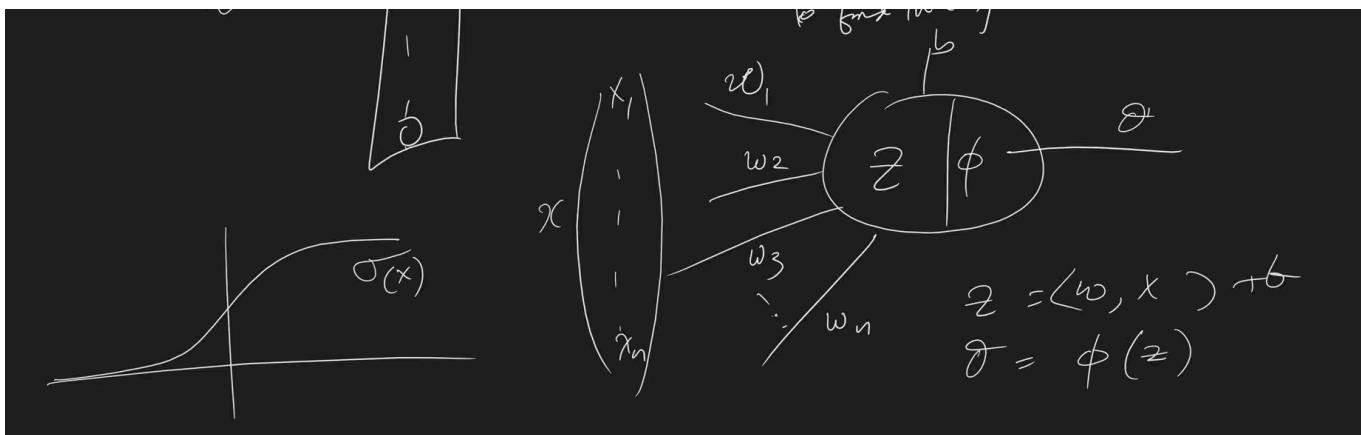
proven to approximate **any** function, they can, theoretically, model **any** causal relationship in the real world. Universal Approximator Theorem guarantees that for any unknown ground truth function  $g$ , a neural network can find an effective function  $f$  that closely approximates it. This  $f$  is found by training the network.

## The Neural Network Training Loop



To find the effective function  $f$ , the network must be trained, which is the process of adjusting its **trainable parameters** ( $\theta$ ), the weights ( $w$ ) and biases ( $b$ ) in its neurons.

## Model Mechanics and Non-Linearity



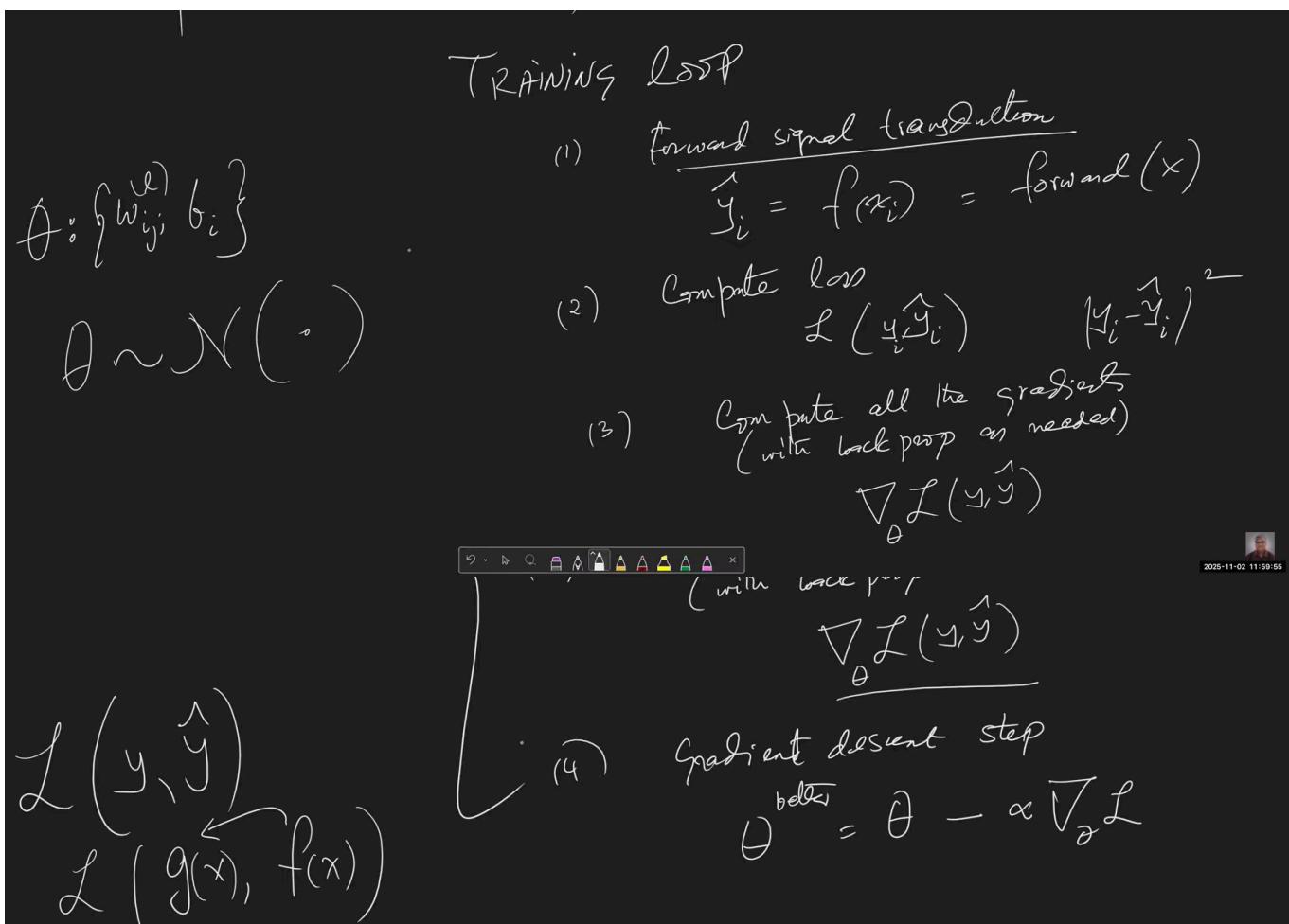
A single neuron's operation is defined by:

$$z = wx + b$$

$$\theta = \phi(z)$$

where  $\phi(\cdot)$  is the activation function. The network gains the power to approximate complex, non-linear functions by combining linear operations ( $wx + b$ ) and a non-linear deformation introduced by  $\phi(\cdot)$ .

## The Four Steps of Model Training



The training process is a systematic, iterative loop that minimizes the divergence between the model's prediction and reality:

1. **Forward Signal Transduction:** The input  $x$  is passed through the network to produce a prediction  $\hat{y}$ , expressed as  $\hat{y} = f(x)$ .
2. **Compute Loss:** The prediction  $\hat{y}$  is compared against the actual ground truth value  $Y$  to calculate the **loss**  $\mathcal{L}$ , which measures the error (e.g.,  $\mathcal{L} = (y_i - \hat{y}_i)^2$  for regression).
3. **Compute Gradients with Backpropagation:** This step calculates the **gradient of the loss** with respect to every parameter:  $\nabla_{\theta}\mathcal{L}(y, \hat{y})$ . The gradient is a measure of a parameter's contribution to the error.
  - **Intuition of Gradient (Slope):** Imagine standing on the ground. If you take a step and your altitude doesn't change, the ground is flat. If your altitude changes, you are on a **slope**. The gradient measures how much height you gain for a single step forward, and since you can walk in any direction, there will be a slope in every direction.
  - **Intuition of Gradient (Sensitivity):** The gradient is the **sensitivity to a step action**. Consider tuning an old analog radio with a **big dial** and a **small dial**. If you're in the middle of static and turn the dial, the signal doesn't change much (low sensitivity). When you start hearing the station, you move in the right direction. The big dial might move you past the station quickly (high sensitivity/bigger slope), while the small dial requires quite a bit of rotation to **gradually inch** your way toward the station (lower sensitivity/smaller slope). The gradient captures this difference in sensitivity for each parameter.
- Gradient Descent Step: The parameters are updated by taking a small step in the opposite direction of the gradient to ensure the loss is minimized.

$$\theta^{\text{new}} = \theta^{\text{current}} - \alpha \nabla_{\theta}\mathcal{L}$$

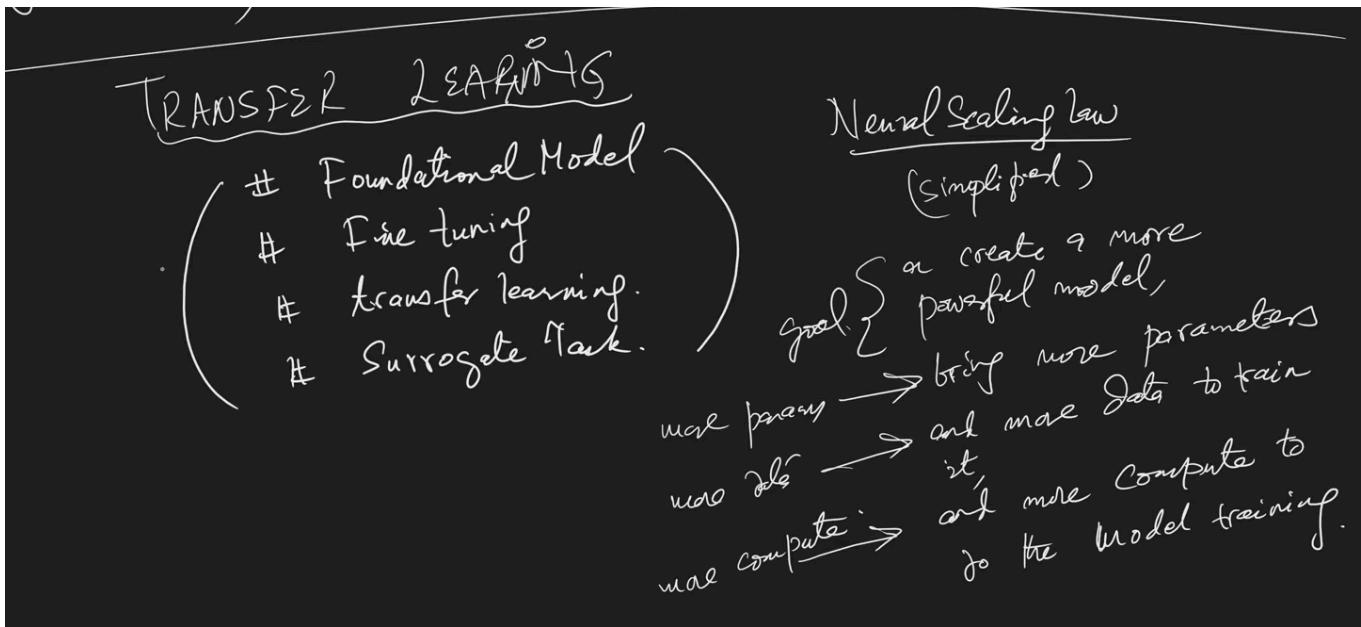
This cooperative adjustment of all parameters causes the function  $F$  to converge towards  $G$ , proving the model to be a useful and effective approximator of the ground truth.

## Null Hypothesis

Model training can only occur if there is a perceived **relationship** between the input ( $x$ ) and the output ( $y$ ). If data points are scattered randomly with no curve or surface to hug (Scenario A vs. Scenario B), then the **Null Hypothesis** is true (no relationship exists). In this case, no matter how large the neural network is, it will **never** be able to predict a relationship because none exists.

---

## Transfer Learning



**Transfer Learning** is the process of utilizing knowledge gained from solving one task and applying it to a different, but related, problem. The concept is analogous to human learning: knowing how to ride a bicycle makes it easier to learn how to ride a motorized bike, even though some **residual learning** (e.g., how to use the engine) is still required for the new task.

The core concepts related to this process are: **Foundational Model**, **Fine-Tuning**, **Transfer Learning**, and **Surrogate Task**.

## The Challenge of Scale and The Neural Scaling Laws

The need for Transfer Learning is driven by the extreme cost of modern model training.

- **Neural Scaling Laws:** A simplified version of these laws dictates that to create a **more powerful model** (with better performance), one must bring **more parameters**, **more data** to train it, and **more compute** to handle the training process.
- **Limitation:** Very few entities (like Google, which trained models such as AlexNet) have access to the necessary data and compute resources. This led to a tradition of organizations training large models and then **releasing them to the public**.
- **The Paradox:** A model is trained with a specific **loss function**, which is the mathematical representation of a particular **engineering goal** (e.g., predicting ice cream sales). The paradox is: how can a model trained for one specific goal be useful for the public, whose goals are entirely different (e.g., a farmer's crop monitoring)?

## The Solution: Foundational Models and Foundational Ability

when you train a model with a foundational ability,  
many downstream tasks look possible as last mile  
problems !!

The utility of a publicly released model lies in the concept of a **Foundational Model**. The key is to train the model on a **Foundational Ability** rather than a specific end-goal. Many seemingly different tasks share a **lowest common denominator** functionality. Such as telling if a **pepper tree has cancer** (for **early detection** against an epidemic) using a drone, classifying a **cloud formation** as "blue sky" versus a **Category 4 hurricane**, looking at a device to tell if it has a **micro-fracture** (fatigue), Identifying a cow, duck, or car. The lowest common denominator is that the model needs to have visual language understanding. For text-based tasks such as determining whether the **sentiment** of a text is positive or negative, or judging whether an **answer** is relevant to a **question** or a *non sequitur*, the downstream model needs to possess natural language understanding.

Domain	Diverse Tasks (Last-Mile Problems)	Foundational Ability (Core Prerequisite)
Text	Sentiment analysis, answering a user prompt, checking for text entailment.	<b>Natural Language Understanding (NLU)</b>
Images	Telling if a pepper tree has disease, classifying a hurricane's category from cloud formations, detecting micro-fractures in metal.	<b>Visual Language Understanding</b> (understanding abstract meaning from pixels)

A **foundational model** is a model trained extensively to master this core prerequisite, the **foundational ability** (e.g., NLU or VLU). This training is extremely difficult and resource-intensive. Once a model possesses this foundational ability, the wide variety of **downstream tasks** become trivial, or "last-mile problems." The model only needs to be shown a few examples (e.g., a happy customer statement and an unhappy one) to quickly adapt its existing understanding to the new, specific task. This is the essence of **transfer learning**: training is shifted from learning the fundamentals to merely adapting them.

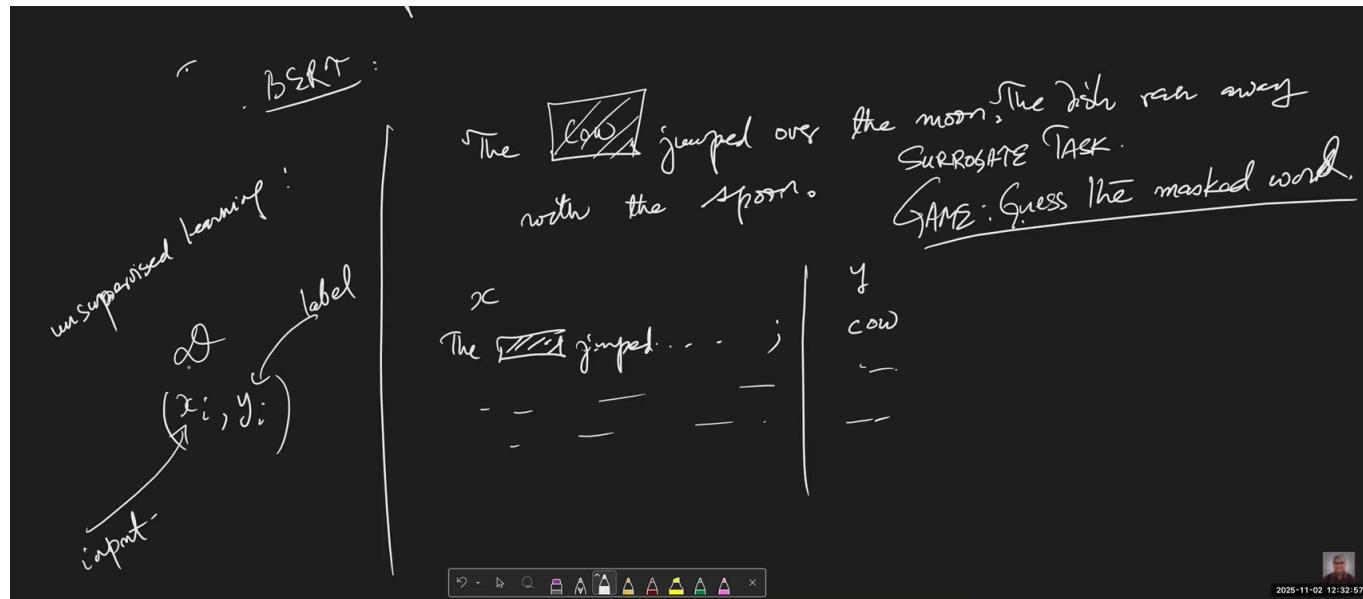
# Training Foundational Abilities through Surrogate Tasks

While mastering a **foundational ability** (like Natural Language Understanding or Visual Understanding) is essential, it is also resource-intensive and difficult. The solution is to use a **surrogate task** during the **pre-training stage** to induce this core ability.

## Surrogate Task: The Forcing Function

A surrogate task is an artificial problem that has **no direct practical utility** on its own, but forces the model to develop the desired foundational understanding as a byproduct of solving it.

### The Example of BERT (Guess the Masked Word )



The data for training language models is obtained by **scraping the web**, which provides a massive collection of raw, unlabelled text. To create usable supervised training data, researchers engage in **data manufacturing** by constructing examples where a word in a sentence is deliberately masked, for instance, turning “The cow jumped over the moon” into “The \_\_\_ jumped over the moon.” Here, the **input** ( $X$ ) is the masked sentence, and the **output label** ( $Y$ ) is the missing word (“cow”). This setup forms a simple **game-like surrogate task**: the model must “guess the masked word.” Although commonly referred to as “unsupervised learning” because the data is not manually labeled, the process is in fact **supervised**, since each input has a clearly defined, artificially created label.

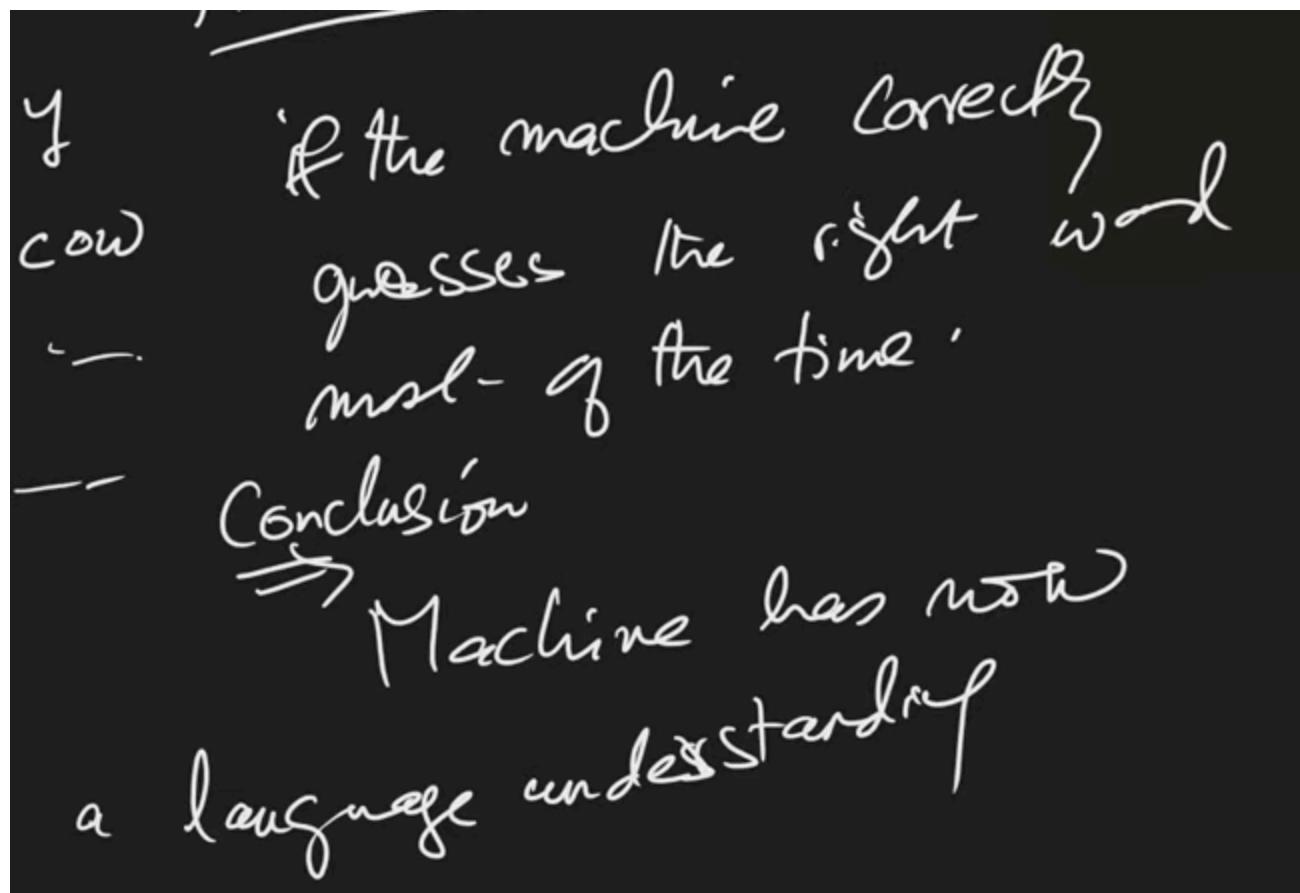
## Inducing Foundational Ability

The value of the surrogate task lies in its effectiveness as a **forcing function** to induce a core ability.

## The Mandarin Child Analogy

Imagine taking an American child and planting them in the Hunan territory of China, playing only the "Guess the Missing Word" game in Mandarin. Initially, the child would struggle. However, if the child begins to consistently guess the correct Mandarin word from a vocabulary of millions, the only logical conclusion is that **the child is getting the hang of Mandarin**, they have achieved language understanding. Statistically, the probability of correctly guessing the right word consistently by pure chance is practically **zero**. The guessing ability is a **measure of language understanding**.

## The Machine's Learning



When a machine plays the game, it produces a **probability distribution** over its entire vocabulary for the masked word:  $P(\text{token} | x)$ . : The training loop uses the loss function to repeatedly nudge the machine to **maximize the probability** (i.e., move the probability mass) around the correct word. The machine does this by constantly changing its **weights and biases** effectively changing the shape of its probability function until it consistently predicts the right word. This ability to make the correct prediction is the practical manifestation of **learning the language's nuances**. **If the machine correctly guesses the right word most of the time, the conclusion is that the machine has now a language understanding.**\*\*\*\*

## Transfer to Last-Mile Problems

A model trained on this foundational ability (e.g., a Natural Language Understanding model), is now ready for **Transfer Learning**. All subsequent tasks (like sentiment analysis) become **last-mile problems**. For instance, the child who understands Mandarin would quickly grasp that the phrase "This fish dish is rotten" indicates a **negative sentiment** with just a few examples of positive and negative reviews. The process of mapping this core foundational ability to your specific domain and specific task is called **fine-tuning** - the final, crucial step in the transfer learning process.

---

## Quantifying Loss and The Enterprise Use Case

The foundation of transfer learning relies on the concepts of the **Foundational Model**, **Foundational Ability**, and the **Surrogate Task**. We can quantify the training of the foundational ability using a specific loss function.

### Loss Function for the Surrogate Task (BERT)

For the surrogate task of "Guess the Masked Word" (like in BERT), the machine aims to predict the probability of the correct word given the context,  $P(\text{cow}|X)$ . The standard loss function used to teach the machine to maximize this probability is the **Negative Log-Likelihood ( $\mathcal{L}$ )**:

$$\mathcal{L} = -\log P(\text{cow}|X)$$

When the machine is **very confident** that the word is correct ( $P \approx 1$ ), then  $\log(1) = 0$ . The loss is zero, meaning the machine made no mistake. If the machine is **not confident** and assigns a very low probability to the correct word (e.g.,  $P = 10^{-6}$ ), the log value is a large negative number. The negative sign in the loss function flips this to a **large positive loss**. By trying to **minimize the loss** during training, the machine is forced to **maximize the probability** of the correct word. This continuous adjustment of weights is the learning process that induces the foundational ability (language understanding).

## Transfer Learning in an Enterprise Context

A common enterprise question is whether a pricing optimization model trained for one industry (e.g., **retail**) can be transferred to a different one (e.g., **software**).

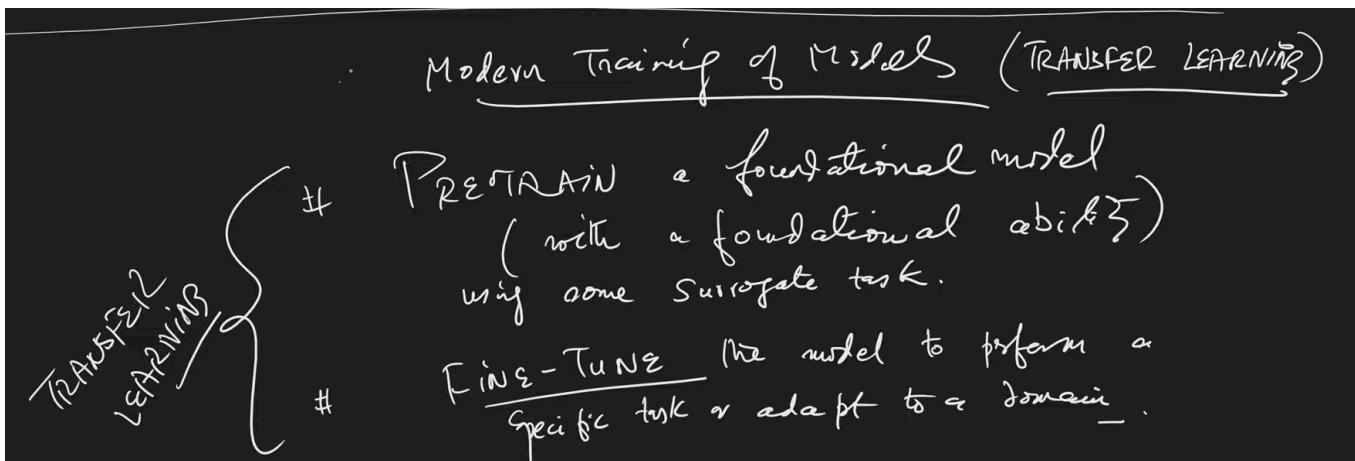
Concept	Explanation	Applicability to Pricing Optimization
<b>Foundational Ability</b>	The core, expensive-to-induce skill that is shared across tasks (e.g., NLU).	In theory, all pricing models share basic economic concepts like <b>price elasticity</b> and <b>demand/supply</b> . This shared understanding is the potential foundational ability.

Concept	Explanation	Applicability to Pricing Optimization
Transfer Learning	Using a pre-trained model as a starting point for a new, specific task.	If the underlying economic "laws" are transferable, you could try using the retail model to start training the software model.
Fine-Tuning	The final step of refining a foundational model on a specific domain's data.	Training the retail model's weights on the specific data of the software industry.

## Conclusion on Pricing Models

In practice, pricing models are often **not highly transferable** because the underlying forces driving consumer decisions and cost structures are vastly different (e.g., **Zara's** quick commerce/fashion vs. **semiconductor** sales driven by sheer productivity/ROI). The weight drift during fine-tuning would be so high that the value gained from the transfer is minimal. Furthermore, the primary criteria for justifying transfer learning (the Neural Scaling Laws: **big compute, big data**) are often absent in traditional regression models like those used for pricing. Training a pricing model from scratch usually takes **barely 20 minutes** on modern systems. **It's faster and more efficient to train a new model from scratch** on the client's data than to attempt complex fine-tuning. Therefore, for such use cases, the typical enterprise solution is to create a general-purpose architecture that trains a separate model on-site for each client's specific data, rather than relying on a transferable foundational model.

## Modern Model Training: The Two-Step Journey of Transfer Learning



The current methodology for training powerful and versatile machine learning models, particularly large language models, is defined by a highly effective **two-step journey** known as

**Transfer Learning.** This approach efficiently separates the creation of general knowledge from its application to specific tasks.

The process begins with **Pre-training**, which involves creating a **foundational model**. This model is endowed with a **foundational ability** (such as a deep understanding of language structure, syntax, and semantics) by being trained on massive, diverse datasets. The training is typically achieved by making the model learn through a **surrogate task**, a large, self-supervised objective that acts as a proxy for general learning, rather than being the final application itself. This initial step is computationally intensive and is often performed by large institutions..

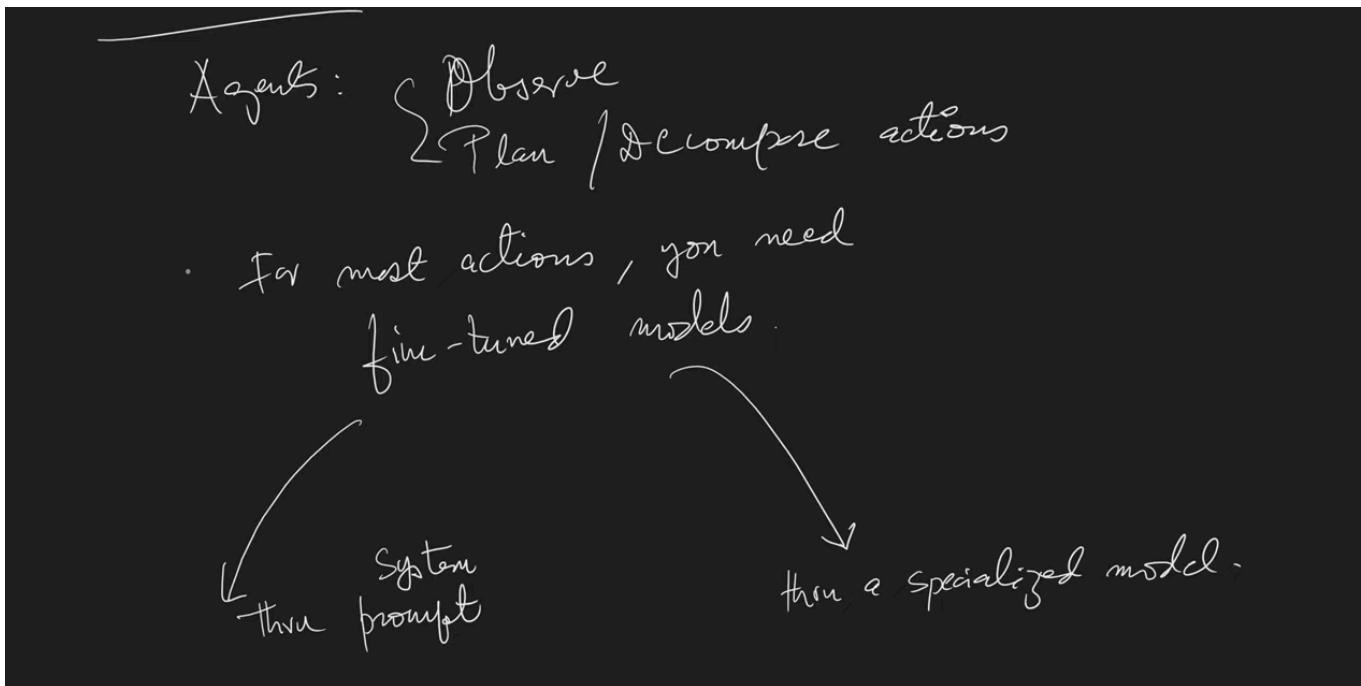
The second step involves adapting the pre-trained model to real-world utility through **Fine-tuning**. **Fine-tuning** takes the foundational model and adjusts its parameters to **perform a specific task** (e.g., sentiment analysis) or **adapt to a specific domain** (e.g., medical texts). This is the stage where individual developers or smaller teams (the "normal peasants") take a model, often downloaded from repositories like Hugging Face and specialize it for their needs. The core idea is **transferring the ability to understand language** (or other foundational knowledge) to a new, specialized function, like analyzing customer reviews.

## **The Unifying Concept: Transfer Learning**

The **technical word for this entire two-step journey**, from foundational knowledge acquisition to specialized application is **Transfer Learning**. It is a name coined for this specific training paradigm where the extensive, initial learning is leveraged and repurposed. **Fine-tuning is a part of Transfer Learning**; it is the specific method used to execute the knowledge transfer. This is why these two steps are conceptually joined together and referred to collectively as the **Transfer Learning approach**.

---

## **The Necessity of Fine-Tuning AI Agents**



The core question driving the evolution of AI applications is **why fine-tune agents**, and the answer lies in moving beyond the current, often unreliable, practice of using simple prompts. Today, most AI agents built using orchestration frameworks like **Crew AI**, **Autogen**, **LangGraph**, or **Google ADK**, are fundamentally just **prompts over an LLM**. This leads to the valid criticism that "so-called agents are nothing but prompts fashionably dressed."

The authentic purpose of an AI agent is to demonstrate robust, reliable, and specialized behavior, following a strict workflow: **observe, plan, decompose (into actions), reason, and act**. In simple (or "baby") projects, a single foundational LLM (like those from OpenAI, Anthropic, or Gemini) can handle the entire task through a prompt. For complex, real-world tasks, when an agent must **decompose** a major goal into individual **granular tasks**, relying on a single frontier model becomes inefficient and prone to error. For each of these individual tasks, we need to have fine-tuned models.

## Why Fine-Tuning Beats Prompt Engineering ?

The central argument for fine-tuning revolves around reliability, deep character integration, and reducing the error rate.

### In-Context Learning vs. Deep Habit

When you write a **system prompt** and give it the user's data (e.g., a proposal for a children's park), both enter the LLM, which then performs **in-context learning**.

Prompt-Based (In-Context Learning)	Fine-Tuned (Ingrained Habit)
<b>First Attempt:</b> The LLM is essentially reading the "script" (the prompt) and playing the role for the first time, improvising as it goes. This is In-Context Learning	<b>Deep Integration:</b> The model's "very basic nature changes" as it has been shown millions of examples of how to behave.
<b>Mediocre Performance:</b> This is analogous to a movie actor being handed the script on the morning of the shoot and expected to act immediately.	<b>Great Performance:</b> This is analogous to a great actor (like Ben Kingsley playing Gandhi) who lives and breathes the role, making the character an inherent part of their nature.
<b>Unreliable Output:</b> The error rate is high because the model is "learning and doing the job for the first time every time the request is invoked."	<b>Low Error Rate:</b> The behavior is baked into the model's parameters through repeated training (supervised fine-tuning or reinforcement learning), correcting mistakes over time.

## Specialized Responsibility

A fine-tuned model becomes deeply imbued with a specific responsibility. - Consider a plan for a children's park with swings and slides. A specialized model is needed to immediately spot key concerns like the environmental impact of the location (next to a lake), regulatory questions, and issues designed to mitigate injuries (e.g., how the main problem with the swing, pillar loosening is addressed). To make an LLM sensitive to this language, you can either write a complex system prompt or bake it into the model. The latter creates a specialized model that keeps a "hawk eye" for these critical details. This deep integration explains why the biggest criticism against today's agentic applications that they can't be trusted because "one moment it produces the right answer, another moment it doesn't", is often valid, as they rely on prompts wearing fancy dresses instead of specialized, fine-tuned models.

## The Last Mile Problem

Fine-tuning models for agents is not overly difficult because it is a **last mile problem**. You are taking a model that already understands human language and simply asking it to specialize (e.g., asking a model to become a legal expert or, humorously, to create a more effective Shashi Tharoor than a prompt ever could).

---

## Frontier Models vs. Fine-Tuning: A Cost-Benefit Analysis

The question of whether fine-tuning is necessary when **frontier models** are becoming smarter and smaller is an essential trade-off between cost, performance, and specialization. While powerful large models can often perform complex tasks via smart prompting (in-context learning), fine-tuning smaller models for core tasks often proves more **cost-effective and reliable** for applications that need to scale.

## The Fine-Tuning Spectrum: Sizing the Brain



The first task in model selection is determining the optimal **size of the brain** (model size) required for the task. This involves a complexity study of the task's **cognitive complexity**.

Typical fine-tunable model size categories range approximately from:

- 2 Billion (2B)
- 7-8 Billion (7-8B)
- 13 Billion (13B)
- 30 Billion (30B)
- 70 Billion (70B)
- 120 Billion (120B)
- 200 Billion (200B)

A practical approach is to **de-escalate** the size: start with a frontier model and a prompt, observe the performance, then scale down to smaller models (e.g., 120B → 70B → 30B), noting where performance suddenly **dips**. The model size before the dip is the sweet spot for fine-tuning.

- **Example Use Cases:**

- **Small Models (e.g., 2B parameters):** Used for quick, first-level filtering of junk content or spam in web searches.
- **Large Models (e.g., 200B+ parameters):** Necessary for tasks requiring complex reasoning, planning, task decomposition, and contradiction/consistency search.

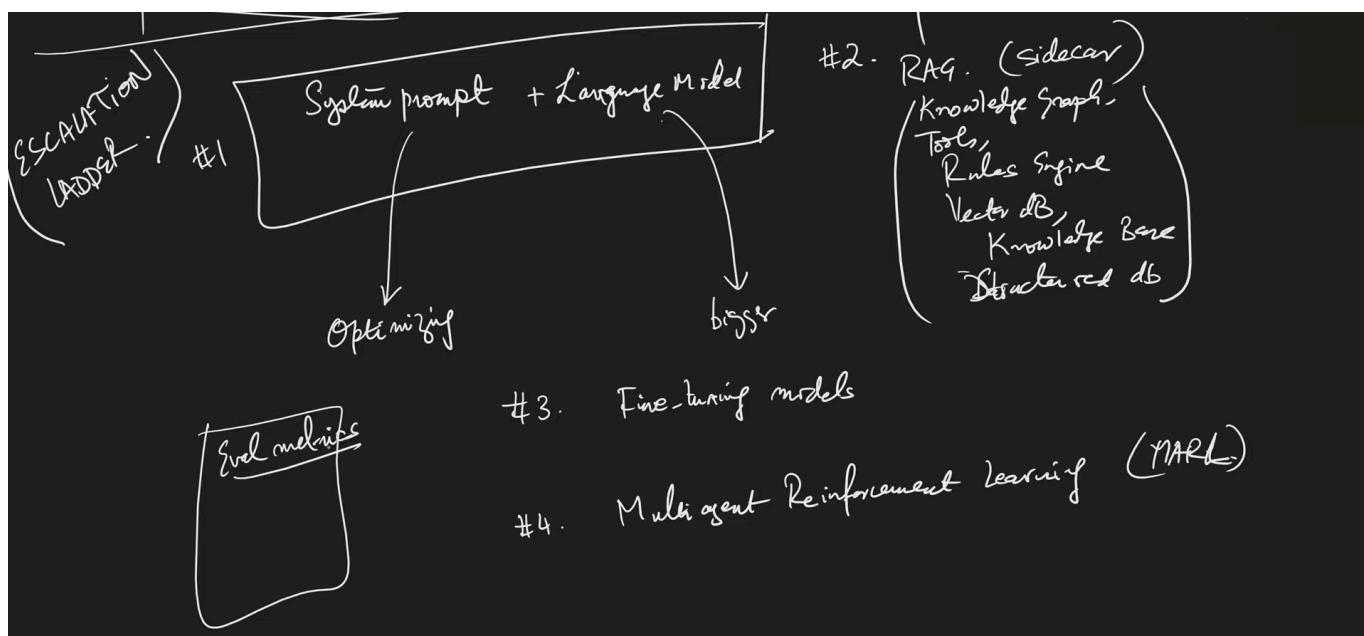
## The Analogy of Specialization

The distinction between using a frontier model with a long prompt and a fine-tuned model mirrors the difference between a generalist and a specialist:

- **Frontier Model + Prompt (Generalist):** This is like interviewing a generalist for a long time to assess their fit for a specific job (e.g., log mining). You have to supply an alarmingly **long set of instructions** (a mile-long prompt) every time, leading to lower reliability.
- **Fine-Tuned Model (Specialist):** This is like hiring a **neurosurgeon**. The interview is short (a few minutes) because their deep, ingrained specialization speaks for itself. The instructions (the system prompt) are **minimal** because the model is already **tuned to do the task**.

As noted, **fine-tuning is akin to using a Digital Signal Processor (DSP)** designed for a specific instruction set, which can outperform a general-purpose CPU (the frontier model) running at a much higher clock speed for that dedicated task. This specialization reduces the **information exchange (entropy)** needed per task invocation.

## The Escalation Ladder for Application Development



When developing an application, one should always use an **escalation ladder** and **stop wherever the evaluation metrics are met**.

1. **System Prompt + LLM:** The cheapest and first step. Optimize the prompt.
  2. **Sidecar Mode (RAG Pipeline):** If (1) is insufficient, bring in sidecars:
    - RAG with **Knowledge Graphs, Tools, Rules Engines, Vector DBs, and Structured Databases.**
    - **Rationale:** RAG is generally a better go-to method for handling **fluid knowledge** because the knowledge base can evolve without necessitating a model re-train.
  3. **Fine-tuning Models:** If (2) is insufficient, fine-tune the core LLM.
    - **Decision Point:** Fine-tuning is preferred when **specialized reasoning or a fundamental change in habit/behavior** is required, not just knowledge retrieval. Models are meant to **reason and think**, not just memorize.
  4. **Multi-Agent Reinforcement Learning (MARL):** The highest and most complex escalation.
- 

## RAG vs. Fine-Tuning: A False Dichotomy

The debate is not truly RAG *versus* fine-tuning; rather, a good RAG system is often built *with* fine-tuning. If the problem is knowledge retrieval, use RAG (the sidecar). If the problem is reasoning, behavior change, or generating precise output, use fine-tuning. While "effortless RAG" (simple, vanilla implementations) is cheap and easy, its quality is often mediocre. To create an effective RAG system, internal components like the semantic embedder and the re-ranker need to be fine-tuned. This adaptation can take as little as 20 minutes but makes a "night and day difference" to the RAG system's quality. This represents the final step where the focus shifts from training one powerful individual model to training a team of agents to play well together—a system-level fine-tuning where multiple agents jointly back-propagate and improve to achieve an overarching goal more effectively.

---

## Distinguishing Fine-Tuning from Knowledge Augmentation

The confusion that **fine-tuning** is simply a combination of a reasoning model plus domain knowledge (achievable with a good prompt and **RAG**) stems from a fundamental misunderstanding of what models should be used for. **Fine-tuning should almost never be used to bake knowledge into a model.**

## Why Fine-Tuning is for Behavior, Not Knowledge

Fine-tuning is a specific technical process in the transfer learning journey: further training an already pre-trained model using a training loop. It is the method for changing a model's *behavior, reasoning, and specialized ability* (like an actor deeply ingraining a character), not its *knowledge base*. Knowledge stored within an LLM's weights is highly compressed, lossy memory. This fact was demonstrated in a landmark case involving ChatGPT. The New York Times sued, claiming the model was plagiarizing their articles by reproducing language belonging to their writers. OpenAI initially argued this was mathematically impossible, stating their model was just weights and biases with no stored essays. The New York Times proved otherwise: they treated the LLM as a black box, fed it the first half of a paragraph, and the model faithfully completed the entire article. This demonstrated that the model was, in fact, storing and retrieving information. The subsequent research confirmed that this knowledge is stored in a compressed, lossy form. When this information is "decompressed" or retrieved, there is a high error rate. It is mathematical illiteracy to use an expensive and unreliable method like fine-tuning to store knowledge when RAG offers absolute preservation of knowledge and reliable retrieval. Those who write articles saying "RAG is dead" because models can be fine-tuned to remember things are fools; they show a lack of reasoning from first principles.

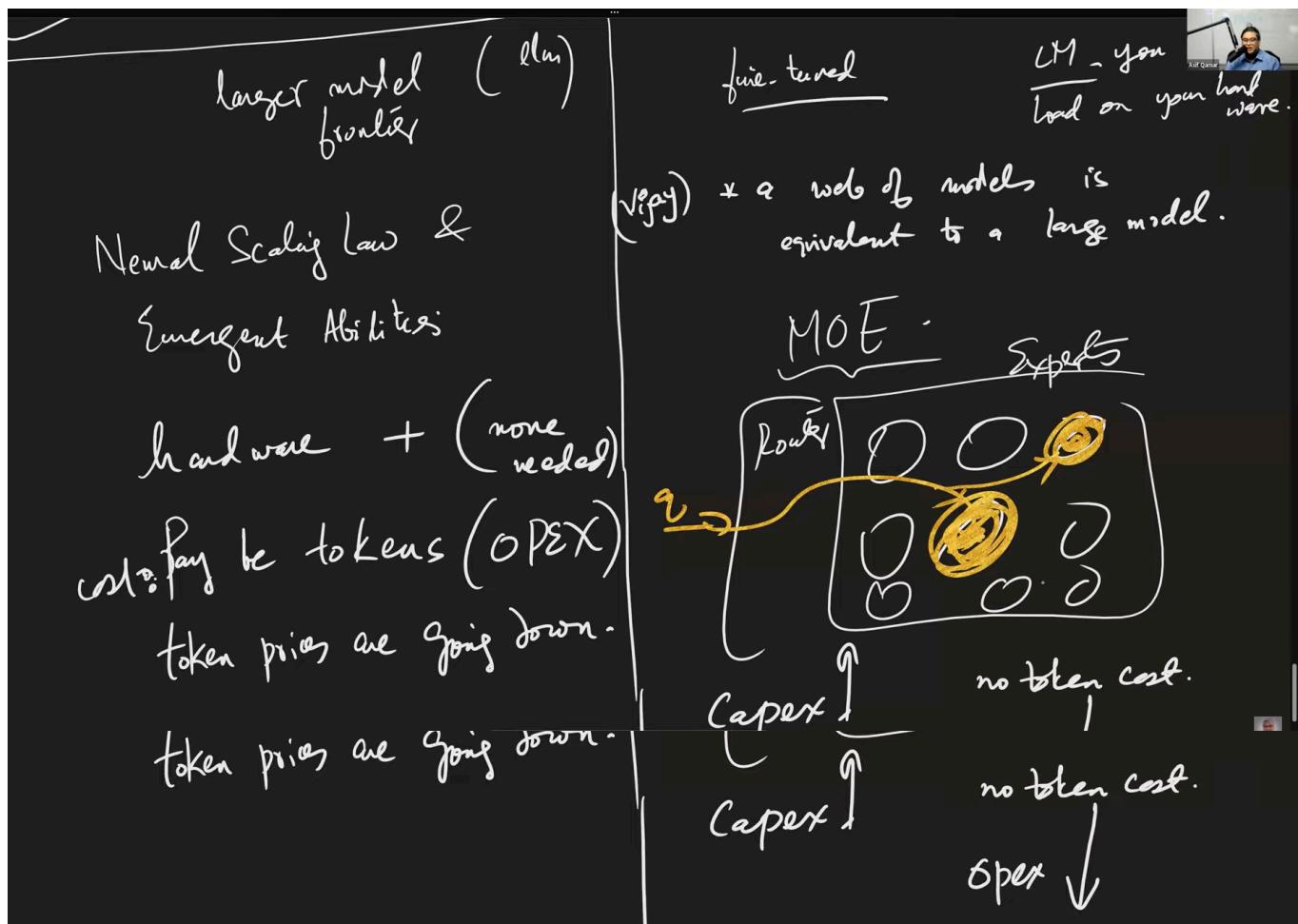
- **Fine-tuning** is like a great musician achieving mastery through practice and corrections. It enhances **conceptual understanding** so the model can apply reasoning.
- **RAG** is like the system you use to practice LeetCode: it's not core understanding, but retrieving an adjacent, already-solved problem from a separate knowledge base to facilitate easier reasoning.

## The Fine-Tuning Stage in the Escalation Ladder

If System Prompt + LLM and RAG are insufficient, then fine-tuning is required. This often means fine-tuning components *within* the RAG system itself, not just the core LLM. Fine-tune individual models within the pipeline, such as the **semantic embedder** or the **re-ranker**, to improve **precise knowledge retrieval**. Fine-tuning an embedder for a specific domain can take as little as 20 minutes but drastically improves RAG quality. Fine-tune specialized LLMs, such as the **query decomposer** (which plans the answer) or the **critique LLM** (which checks for hallucination).

---

## LLMs vs. Fine-Tuned Models



The debate pits the philosophy of "**Larger is Better**" (relying on scale and emergent abilities) against the strategy of "**Specialization and Distribution**" (relying on fine-tuning and smaller, focused models).

## Arguments for Large Frontier Models (LLMs)

This camp argues for using powerful, pre-trained models, often accessed via API (paying by token).

Concept	Argument in Favor	Key Evidence/Analogy
<b>Neural Scaling Laws &amp; Emergent Abilities</b>	As models increase in size, data, and compute, they develop <b>emergent abilities</b> —intellectual leaps that aren't present in smaller models.	<b>Analogy:</b> When facing a hard math problem, you'd rather hire a <b>smarter mathematician</b> than spend years teaching a weaker one.
<b>Domain Mastery</b>	Sufficiently large models can overcome domain-specific dialect	<b>Bloomberg LLM Case:</b> Bloomberg trained a domain-specific LLM on financial data. While it beat earlier models like GPT-3.5, <b>GPT-4</b> (a

Concept	Argument in Favor	Key Evidence/Analogy
	issues without explicit retraining.	much larger frontier model) surpassed it in every financial benchmark, proving that <b>greater scale can lead to superior domain understanding.</b>
<b>Cost Structure</b>	Eliminates large upfront hardware costs (CAPEX ).	Users <b>pay by the token</b> (OPEX). Companies assure that token prices continually drop year-over-year. <b>No dedicated hardware</b> is needed for inference.
<b>Simplicity of Architecture</b>	You can solve complex problems with one model and an iterative loop (like the <b>ReAct</b> loop), which is difficult for smaller models to manage.	The goal is a simpler architecture: one large model can handle <b>reasoning, planning, and decomposition</b> without needing multiple specialized agents.

## Arguments for Fine-Tuned Models (SLMs / Web of Models)

This camp argues for using smaller, specialized models that are trained or fine-tuned for specific tasks within an agentic system.

Concept	Argument in Favor	Key Evidence/Analogy
<b>Distributed Intelligence</b>	A collection of specialized, fine-tuned models can collectively achieve or exceed the abilities of a single frontier model.	<b>Insight:</b> A web of fine-tuned models is conceptually equivalent to a single large model. The entire philosophy of <b>AI Agents</b> is built on this insight: <b>distributed AI</b> where $2 + 2 > 4$ .
<b>Efficiency and Cost at Scale</b>	Although fine-tuning requires upfront hardware investment or rental (CAPEX), the operational cost for inference is drastically lower.	<b>Cost &amp; OPEX:</b> After the initial training, the operational expense is near <b>zero</b> (just electricity/cooling). There is <b>no recurring token cost</b> , making it far more economical for high-volume, repetitive tasks.
<b>Exploiting MoE Architecture</b>	Modern large models are often <b>Mixture of Experts (MoE)</b> models, where only a small sub-model (expert, e.g., a	<b>The Specialization Trap:</b> If an agent is designed for one specific task, the large model will repeatedly activate the <b>same one or two internal experts</b> . You're essentially paying for the entire 600B model's

Concept	Argument in Favor	Key Evidence/Analogy
	30B parameter expert within a 600B model) is activated per query.	infrastructure just to run a 30B expert. It's more efficient to <b>take a 30B model in-house and fine-tune it to be the ultimate guru</b> on that topic.
<b>Specialization &amp; Reliability</b>	Fine-tuning a model on a narrow use case (e.g., an agent's specific role) makes it more reliable and accurate for that role.	<b>Analogy:</b> Fine-tuning is like taking a capable actor (Ben Kingsley) and dedicating time and training to make him fully embody a specialized role (Gandhi), ensuring peak performance.

## The Crucial Trade-Offs

The choice between the two fundamentally boils down to a few key trade-offs. The primary difference lies in cost structure: large models involve OPEX (recurring operational expenses such as token usage), whereas fine-tuned models require CAPEX (upfront investment in hardware or rental). In terms of flexibility and generality, large models are highly adaptable and capable of performing a wide range of tasks, while fine-tuned models are more specialized, excelling at a single, well-defined task. When considering latency and speed, large models may experience higher latency due to their size and dependence on API calls, whereas fine-tuned models tend to run faster and can operate on smaller, optimized hardware. Finally, the difficulty of implementation differs significantly, large models are easier to use through plug-and-play APIs, while fine-tuned models demand more effort and skilled engineers for proper setup and optimization.

## A2A Conceptual Deep Dive

### The Distinction: Client vs. Agent

A key discussion clarified the difference between the test client and an agent.

- The client (`A2A_SimpleClient`) is **not an agent**.
- It is analogous to a **REST client** used to test a REST API. It is a simple program for sending messages to the actual agents.
- The client *uses* the A2A protocol, but this complexity is encapsulated within the ADK library.

#### Note

- **Analogy:** In a real-world application, a user-facing chatbot might be the interface. When a user types into the chatbot, the chatbot's backend service (e.g., a Python web

service) would act as the "client," taking the user's message and calling the A2A agents.

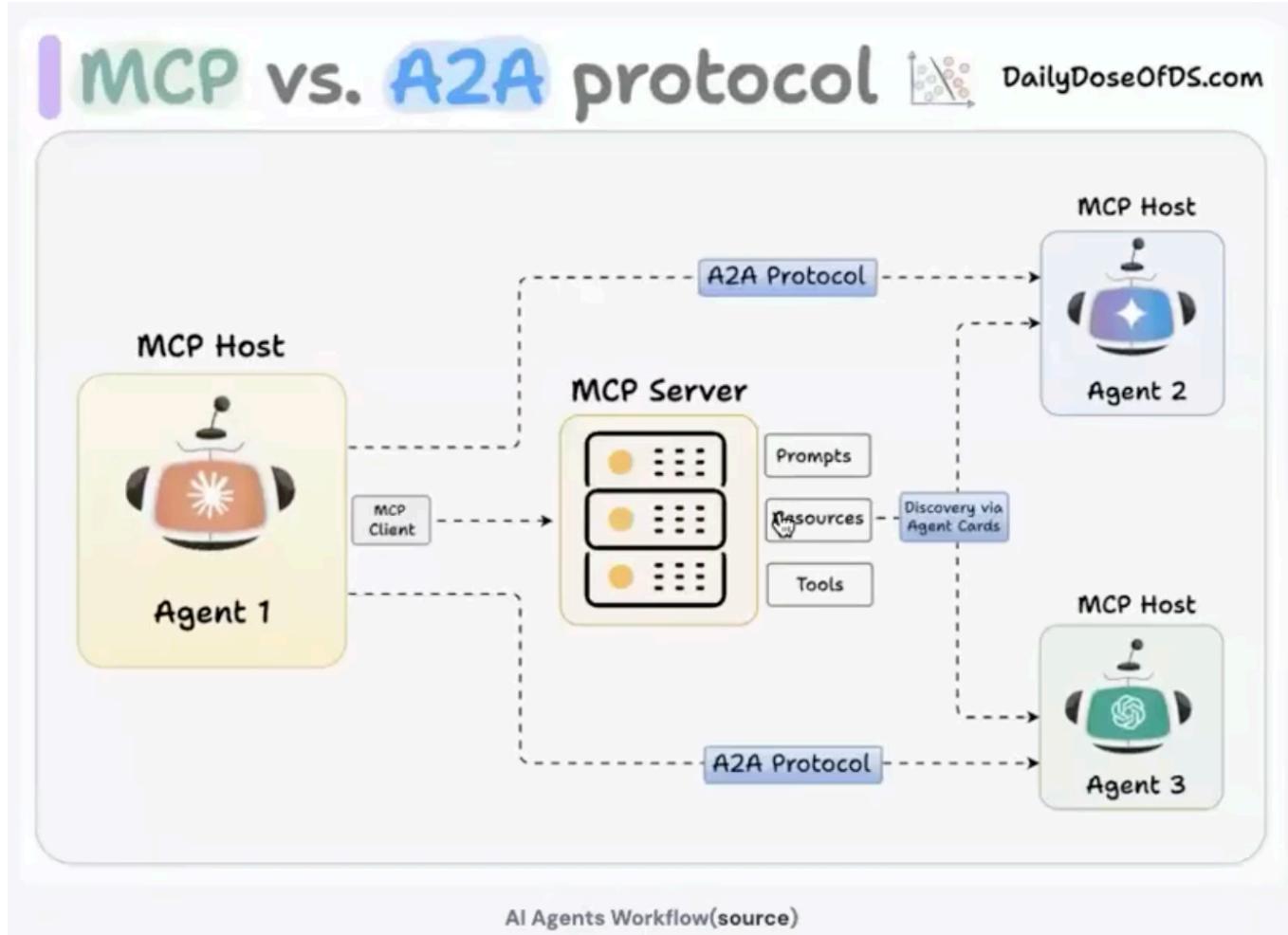
## A2A vs. Single-Process Frameworks (CrewAI, LangGraph)

The primary philosophical difference between A2A and frameworks like CrewAI or LangGraph is **scalability and distribution**.

- CrewAI and LangGraph are highly effective for creating agents that operate within a **single system** or a **single process space**.
- A2A is **fundamentally different** because it is designed from the ground up for a **distributed architecture**.
- This distributed nature allows for true multi-agent communication **across enterprises**. For example, an agent running at one company (e.g., Facebook) could securely communicate with an agent running at another company, each hosted on different IPs and ports.

## Interoperability: A2A, ADK, and MCP

A2A and the ADK provide a higher-level abstraction that can integrate with other systems.



- The A2A protocol can be used to communicate with agents built in *other* frameworks, such as LangGraph. The protocol acts as a translation layer.
- As shown in the architecture diagram, A2A agents can also communicate with and consume tools exposed via **MCP (Model Context Protocol) servers**. An agent can communicate with another *agent* (using A2A) and also call a *tool* (using MCP).
- A single host system could simultaneously host agents (using A2A) and expose tools (as an MCP server).

## Agent Discovery: Peer-to-Peer vs. Centralized

A critical challenge in a distributed system is discovery: how an agent finds other agents. Two models were discussed:

- **Peer-to-Peer (P2P) Discovery:** This is the model used in the demo and supported by the current A2A design. The agent must know the address (IP/hostname) of the agent it wants to talk to in order to request its Agent Card.
- **Centralized Discovery:** This is the long-term vision for a true "internet of agents".
  - **DNS Analogy:** A centralized broker or registry would function like **DNS** for agents.

- An agent could query this central service (e.g., "find me all agents with 'health service' capabilities") to discover agents it did not previously know about.
- **API Gateway Analogy:** This concept is very similar to an API Gateway where services are registered to be discovered.

## The Future of A2A: Standardization and Security

The multi-agent space is currently experiencing "protocol wars," with different companies like Google (A2A), IBM (ACP), and Cisco all proposing standards.

- **Analogy (CORBA vs. JSON):** This is similar to past competitions like CORBA vs. DCOM, where the simpler protocol (JSON/REST) ultimately won.
  - **Open Source Convergence:** In a positive development, these competing protocols are being handed over to open-source bodies like the **Linux Foundation**. Google's A2A is now an open-source project. This convergence will help establish a single, mature standard.
  - **Security:** Security and authentication are the "hottest topics" currently being solved. When agents from different companies interact, they must authenticate and authorize each other, likely by integrating with existing enterprise systems like Okta.
- 

## Advanced Agentic Concepts

### New Framework: Microsoft's Agent-FL for Fine-Tuning

A new Microsoft framework (released early in the week) was introduced, enabling reinforcement learning-based training for LLM-based agents.

- This framework addresses the concept of **agent fine-tuning**, which is distinct from simply fine-tuning a model for knowledge.
- The goal is to fine-tune an LLM to add a *skill*, such as the ability to perform **planning** or **reasoning**.
- The framework appears to work by capturing traces and logs from existing agents (built in LangChain, AutoGen, etc.) and using those traces to train the LLM, making the agent fine-tuning process simpler.

## The Foundation of Tool Calling

A foundational concept for building agents is **tool calling** (or function calling). This is a protocol that allows a client application to tell an LLM about the functions (tools) it has available in its own local space. The LLM can then instruct the client application on *which* function to call to get an answer, which the client then executes and sends back to the LLM. Understanding this

protocol at a deep level (e.g., by tracing the "on the wire" communication) is essential for building robust agentic systems.

## Agent Design Pattern: Planning and Execution

A discussion on agent design patterns highlighted the importance of moving beyond a simple ReAct (Reason-Act) loop.

- A **Planning and Execution** pattern introduces a dedicated **planning step**.
- Instead of immediately trying to act, the agent first calls the LLM with a prompt to "come up with a simple step-by-step plan".
- This plan is then used to guide the execution agent (e.g., a ReAct agent). This explicit planning is a more robust approach.

## Agent Design Pattern: Re-planning

Re-planning is a more advanced pattern that addresses failures in the initial plan.

- **Problem:** A standard ReAct agent can get stuck in an infinite loop or fail after hitting a `max_retries` limit. At this point, the process simply errors out.
- **Solution:** The re-planning pattern adds a node to the graph that catches this failure. Instead of stopping, it triggers a **re-plan** (e.g., by calling the planner LLM again with the new context of the failure), and the process starts over with a new plan. This creates a more resilient, self-correcting loop.

## Agent Design Pattern: Reflection

Reflection is a pattern where an agent assesses the *quality* of its own work.

- This involves prompting an LLM to "observe its past steps" and the results from its tools.
- The agent then "reflects" on these outputs to determine if they are high-quality or if they meet certain criteria (e.g., user preferences).
- If the output is low-quality, the agent can loop back to get a better answer, similar to re-planning but focused on result quality rather than execution failure.

## Example: The "Smart Home" Reasoning Project

An upcoming project about a "smart home" was used to illustrate the practical application of these reasoning patterns.

- **The Problem:** Current smart home alerts are based on rigid `if-then-else` logic (e.g., IF `back_door_open` THEN `send_alert` ).
- **The Goal:** Build a system that uses **reasoning** to provide more intelligent alerts.

### Note

- **Analogy (Sophie the Dog):** A user might want an alert if their dog, Sophie, is stressed. However, they *don't* want an alert if the back door is open *while* the main gate is closed, because they know this is a safe scenario. A simple `if-then-else` system cannot handle this complex, context-dependent logic, but a reasoning-based agent can. This highlights the trade-off: `if-then-else` is cheap, whereas LLM-based reasoning incurs token costs but provides far more flexible and powerful capabilities.

## Practical Career and Skill Discussion

### Real-World Adoption of Complex Patterns

A practical question was raised: Are companies *really* implementing these sophisticated, multi-step agentic patterns (like re-planning and reflection) today?

- The answer is largely **no**.
- The majority of enterprise work (70%+) currently consists of Proof-of-Concept (PoC) projects, which often use simple, 10-line RAG solutions.
- True, complex agentic systems are currently being built primarily by companies whose core business *is* AI (e.g., Perplexity, Cursor).
- For most companies, these advanced patterns will become more relevant as their systems mature and scale to thousands of users, creating a need for optimization and reliability.

### Key Skill Sets for the Evolving AI Field

The discussion identified several key skill sets for individuals navigating the "AI transformation":

- **Data Preparation:** The skill of preparing data for fine-tuning remains critical. This involves understanding data fabric, data lakes, and how to structure data to build a "data factory" for AI.
- **RAG (Basic and Advanced):** RAG is a core component. This splits into basic RAG (which can be learned quickly) and advanced RAG (which may take six months to master).
- **Fine-Tuning:** This is emerging as a specialized, separate skill set. Teams are being formed specifically to handle the fine-tuning of models (especially smaller LLMs).
- **Agentic Skills:** This refers to the ability to build the co-pilots and backend systems, integrating RAG, agents, and knowledge management.

The industry is currently democratizing AI, moving from generative AI (just creating content) to **reasoning AI** (systems that can plan and solve problems).

---

## The Case Against Fine-Tuning for Knowledge

A foundational principle is the separation of *knowledge* and *ability* in AI systems. A strong directive is to **not** use fine-tuning to bake knowledge into a model. Instead, knowledge should be externalized and managed within Retrieval Augmented Generation (RAG) systems.

This does not mean fine-tuning is useless; on the contrary, all high-quality, "good" RAG systems utilize fine-tuned models. Mediocre RAG systems are those that use generic, off-the-shelf downloaded models. The reasons for not embedding knowledge directly into model parameters are numerous and grounded in engineering, thermodynamics, and information theory.

## Knowledge Obsolescence and Cost

Knowledge is fast-evolving and subject to obsolescence; what is true yesterday may not be true today. In a real enterprise, tuning models overnight to reflect these changes is not feasible. The process is expensive; tuning a large model, even on an 8-GPU cluster, can easily require three to four days.

## Risks: Catastrophic Forgetting and Mode Collapse

The fine-tuning process is not risk-free. It introduces significant dangers, primarily **catastrophic forgetting** (where the model loses previously learned abilities) and **mode collapse**.

### Note

Analogy: Neurosurgery

Fine-tuning is aptly compared to neurosurgery. Neurosurgeons are paid exceptionally well because every action they take is irreversible and carries the risk of profound damage.

Fine-tuning experts face the same high-stakes, high-risk environment, which is why their compensation is often comparable.

## Talent and Practicality (The Escalation Ladder)

Proper fine-tuning requires specialized talent that may not be available. In contrast, virtually anyone can build a mediocre RAG system using downloaded models, which may be sufficient to meet initial business goals. This introduces the "escalation ladder" concept: begin with the simplest solution and only escalate in complexity as needed.

# Thermodynamic Disconnect: Knowledge Entropy vs. Model Capacity

Knowledge inherently and exponentially expands. At a civilizational level, it is projected that **more knowledge will be created in the next two years than in the entire history of mankind.**

Models, however, have a finite capacity defined by their parameters (e.g., 70B, 400B, 600B). This creates a fundamental disconnect: knowledge has an expanding entropy, while models have a finite capacity.

## Note

Analogy: Pouring the Ocean into a Bucket

Attempting to bake all expanding knowledge into a finite model is like trying to pour the entire ocean into a bucket. The bucket will inevitably overflow. From a thermodynamic and information-theoretic perspective, this approach is "insanity".

## Parameter Limits and Grokking (Lossy Compression)

Recent research indicates that any single parameter in an LLM has a memory capacity of approximately **3.6 bits**. Once this limit is reached, the LLM stops *memorizing* and begins to *generalize*.

This phase transition is known as "**grokking**". When an LLM groks, it learns patterns and generalizations but *loses specifics*. Therefore, fine-tuning for knowledge is inherently a **lossy compression** process.

## Note

Analogy: The Lazy Professor and Exam Memorization

This is similar to students who dislike a subject and have a lazy professor who repeats exam questions. The students will create a question bank, memorize the perfect answers, and subsequently forget the subject, having never truly learned (or grokked) it. LLMs initially do this (memorize), but when overwhelmed, they are forced to grok (generalize).

## Data Provenance and Legal Requirements

Because fine-tuning is a lossy compression, the exact, original information cannot be perfectly decoded back. From a legal and governance perspective, this is a disaster.

RAG systems, by contrast, can retrieve a specific fact and provide **provenance** or **evidence** to substantiate an answer. An LLM relying on its internal parameters *cannot* do this; when asked for evidence, it will hallucinate. Frontier models use RAG as a backend precisely to prevent this and provide grounding.

---

## The Fundamental Dichotomy and System Architecture

This leads to the central architectural principle:

- **The Fundamental Dichotomy:** LLMs should be fine-tuned to develop **abilities** (e.g., domain-specific reasoning, task-specific skills). **Knowledge** should be externalized and poured into **RAG systems**. LLMs are *not* knowledge containers.

## What is Fine-Tuning?

Fine-tuning is the process of taking a pre-trained model through the *further training process*. It involves the complete four-step training loop:

- Forward signal transduction.
- Loss computation.
- Loss gradients computation (sensitivity of the loss respective to each parameter).
- Parameter update via gradient descent.

## The Maturity Escalation Ladder

A mature AI system is built by progressively escalating its capabilities.

- **Level 1: System Prompt.** Attempt to solve the task using only a system prompt. This is risky for mission-critical tasks due to the stochasticity of the model.
  - **Analogy: The Last-Minute Actor.** Giving an actor their script on the day of the performance will not yield a good or reliable job.
  - **Murphy's Law** dictates that this approach will fail precisely when the most important customer is present.
- **Level 2: RAG.** Back up the prompt with a RAG system (implied).
- **Level 3: Systemic Fine-Tuning.** Begin fine-tuning *all* components of the system: the LLM, the embedding model, the reranker, the semantic cache, etc.. AI is about **system architecture**, not a single model. Apply **Amdahl's Law**: identify the greatest weakness (slowest link) and fine-tune that piece to get the biggest "bang for the buck".

- **Level 4: Reinforcement Learning.** For highly complex situations, the final escalation is to use multi-agent reinforcement learning.
- 

## Debate: Frontier Models vs. Fine-Tuned Local Models

A vigorous debate exists on whether to use large frontier models or smaller, locally fine-tuned models.

### Arguments for Frontier Models (e.g., GPT-4):

- **Performance:** They leverage neural scaling laws and exhibit powerful emergent behaviors.
- **Example:** The Bloomberg financial LM was a significant effort, but its utility was challenged by the release of GPT-4.
- **Cost (CapEx):** There is almost zero **CapEx** (capital expenditure). No hardware purchase is needed.
- **Cost (OpEx):** You pay by the token, and token prices consistently decrease, lowering implementation costs over time.

### Arguments for Fine-Tuned Local Models:

- **Cost:** While **CapEx** is higher (you must buy or rent tuning hardware), the **OpEx** (operational expenditure) is lower once the model is running on your own inference system.
  - **Architecture:** Frontier models are often internally a Mixture of Experts (MoE) anyway. You may *think* you are invoking a 671B model, but you are often invoking smaller, specialized models internally (e.g., two 32B models). Their architecture should guide yours.
  - **Sum of Parts:** Agentic systems are founded on the idea that distributed intelligence is *more than the sum of its parts*.
  - **Analogy: Physicist vs. Biologist:** A fine-tuned "physicist" agent *is* a physicist because it was trained on physics data; a "biologist" agent *is* a biologist.
  - **Governance:** This specialized approach provides **provability, verifiability, traceability, observability, and governance**. You get *none* of these when using a single, large, black-box model to do all jobs.
- 

## Project: The Patient Advocacy Agent

The project is to create an agentic system that functions as a "Patient Advocate".

# Motivation and Narrative

The target audience is frontier villages in the Global South, where the physician-to-patient ratio can be worse than 1:1000 or even 1:10,000.

- Story: Primary Care Units

In many such areas, the government establishes primary care units staffed with employees (e.g., a compounder, nurse, physician). However, the defining quality of these units is that "government opens them, employees close the door". The staff disappears, drawing free salaries, leaving the facility locked.

The project aims to create an accessible unit (a camera, microphone, and speaker) to provide sensible dermatological guidance. This project is designed to be one of social utility. This aligns with the objective from MIT, whose motto ("Ma's moto") is "**technology and the service of mankind**". The goal is to build a project of actual, tangible value that could be deployed to make an impact.

## Core Rule: Do Not Play Doctor

This is the most critical rule. The system is *not* a doctor. The unauthorized practice of medicine is illegal. The agent's role is to provide guidance *where none currently exists* and act as a "patient advocate," similar to a nurse's assistant who remotely prepares a case history.

## Medical Framework: SOAP

The agent will follow the medical **SOAP** (Subjective, Objective, Assessment, Plan) framework.

- **S (Subjective):** The patient's narrative, captured via voice. The agent must filter medically relevant information (e.g., "itching") from irrelevant noise (e.g., "my child was crying").
- **O (Objective):** The medically relevant information, including a picture of the skin condition.
- **A (Assessment):** An evaluation of what the condition *could be*, presented as a set of possible **ICD codes** with supporting evidence.
- **P (Plan):** Contextual guidance for the patient in a remote village. This *must* include transportation directions, urgency, and what to expect.

## Required Outputs

The agent's job is to create a complete case history, delivered in two versions:

- **For the Physician:** A formal case history (including the stored image) sent to a remote healthcare facility.

- **For the Patient:** An explanation in simple, layman's language (via voice) detailing the options and plan.

## Interface and Technical Requirements

- **Interface:** The patient interface **must be voice-only**. No text typing.
- **Rationale:** The target user base is presumed to be illiterate. Illiteracy is no longer a handicap in a world with microphones and cameras.
- **Language:** Must support a minimum of 5 languages, including language detection. The ambitious goal is 100 languages.
- **Privacy:** The agent must explicitly ask for and receive **permission** before taking a picture ("permission-gated pictures").
- **Standard:** **Do not use a Streamlit interface.** The goal is a *real product* deployed on the web.
- **Technology:** This is a classic agentic RAG project requiring multimodal search.
  - Use **WebRTC** for the camera.
  - Use the **SCIN** database from Harvard University ( $\approx$  2GB of text and images) for RAG.
  - Use models like **MedGemma** and **SigLIP-2**.
- Advanced Methodology (For students from previous workshops):

For students who have previously taken the semantic search and language model workshop, higher expectations apply. These students are expected to fine-tune the multimodal embedding model for the retrieval system.

- **Training Method:** The training should be conducted using a **contrastive loss** function.
- **Goal (Isotropy):** The objective is to achieve a more **isotropic distribution** of the embeddings within the unit hypersphere. This separation of concepts in the vector space leads to better and more reliable retrieval results.
- **Testing Method (Clustering):** A tangible method to test the correctness of the embedding model is to run clustering on the resulting image embeddings. A successful, well-trained embedding model will demonstrate that images of the same diagnosis (e.g., **all images of eczema**) **naturally cluster together** in the embedding hyperspace.

## Scope and Common Sense

The system must balance guidance with the "do not play doctor" rule.

- **Critical Cases:** If a malignancy is suspected, the agent must urge immediate action.
- **Non-Critical Cases (De-escalation):** The system must have a list of critical conditions. If the issue does *not* match, it can de-escalate.

- **Example (Blue Paint):** If it looks like paint, the agent should interrogate: "Did you try washing your hand?".
  - **Example (Tattoo):** The agent must be smart enough to ask, "Is that a tattoo?".
  - **Example (Mild Acne):** Reassuring a teenager that it is just acne is a valuable service. Suggesting "take a bath with Lifebuoy soap" is acceptable local advice, not a medical prescription.
  - The agent's role is to *do the leg work* for the patient: gather information, connect them to a provider, and arrange logistics.
- 

## Key Insights

- **Goal of Model Training:** Training a model means finding an **effective function**,  $f(x)$ , to **approximate** an unknown underlying truth,  $g(x)$ . The goal is to make  $f \approx g$ , creating a useful representation of reality from observed data.
- **The Approximation Concept:** All models are considered technically "wrong" because they are only **approximations**, not the true ground truth function. However, they are **useful** if they approximate the underlying reality closely enough for practical application.
- **Regression Example (Ice Cream):** In predicting ice cream sales based on temperature, the model finds  $f(x)$  to predict sales, despite **unknown unknowns** (e.g., wind, holidays). Model quality is measured by metrics like **Mean Squared Error (MSE)** and **Mean Absolute Error (MAE)**.
- **Classification Example (Cows/Ducks):** Classification involves drawing a **decision boundary** to separate data points (like cows from ducks) based on features. The error rate for these models is measured using metrics such as **accuracy**, **precision**, and **F1 score**.
- **The No Free Lunch Theorem:** This principle asserts that **no single algorithm** is universally superior to all others across all domains and use cases. The **effectiveness** of an algorithm is specific to its application and problem type.
- **Universal Approximator Theorem:** A standard neural network with a **single hidden layer** can effectively approximate **any continuous function** to an arbitrary degree of accuracy. This theorem guarantees that neural networks can theoretically model **any causal relationship** in the real world, as causality is expressed through underlying functions.
- **Neural Network Training Loop:** Training involves a systematic, iterative process to find the effective function  $f$  by adjusting the network's **trainable parameters** ( $\theta$ ), the weights ( $w$ ) and biases ( $b$ ). The network gains power to model complex, **non-linear** functions by combining linear operations with a non-linear **activation function** ( $\phi(\cdot)$ ).
- **The Four Training Steps:** The loop starts with a **Forward Pass** to produce a prediction  $\hat{y}$ , followed by **Compute Loss** by comparing  $\hat{y}$  to the ground truth  $Y$ . **Backpropagation** then calculates the **gradient of the loss** ( $\nabla_{\theta}\mathcal{L}$ ), which measures each parameter's contribution

to the error. Finally, **Gradient Descent** updates the parameters by stepping in the opposite direction of the gradient to minimize loss.

- **Gradient Intuition (Sensitivity):** The gradient represents the **sensitivity** of the loss to a small change in a parameter, indicating how much a parameter contributes to the error. The parameters are updated in the direction that **minimizes** this error (loss).
- **Transfer Learning Necessity:** Due to the extreme costs and resource requirements dictated by **Neural Scaling Laws** (more parameters, data, and compute for better models), organizations release large models to the public. This process is driven by the need to reuse and adapt existing, resource-intensive models.
- **Foundational Models and Ability:** The public utility of these large, pre-trained models lies in their mastery of a **Foundational Ability** (e.g., Natural Language Understanding or Visual Language Understanding). Once mastered, diverse, specific downstream tasks become trivial "**last-mile problems**," requiring only **fine-tuning** rather than training from scratch.
- **Surrogate Task (Forcing Function):** A surrogate task is an **artificial problem** used during the resource-intensive **pre-training stage** that has no practical utility on its own. It acts as a **forcing function** to *induce* a desired **foundational ability** (like language understanding) as a necessary byproduct of solving it.
- **Example: Masked Word Prediction:** In language models like BERT, data manufacturing turns raw text into supervised examples by masking a word (e.g., "The \_\_\_ jumped over the moon"). The machine must **maximize the probability** of guessing the correct word, which it does by constantly adjusting its parameters; this ability is the **measure of language understanding**.
- **Loss Function for Foundational Ability:** The standard way to quantify learning for the "Guess the Masked Word" task is using the **Negative Log-Likelihood** ( $\mathcal{L} = -\log P(\text{cow}|X)$ ). Minimizing this loss forces the machine to **maximize the probability** ( $P$ ) it assigns to the correct word, indicating high confidence and successful learning.
- **The Two-Step Journey (Transfer Learning):** The modern, efficient methodology for training large models is a **two-step journey** called **Transfer Learning**. Step one is **Pre-training** to create a **Foundational Model** with a foundational ability (via a surrogate task). Step two is **Fine-tuning**, where the model is adapted for a **specific task or domain**.
- **Enterprise Transferability (Pricing):** While concepts like **price elasticity** are theoretically transferable, models like retail pricing are often **not highly transferable** to domains like software. Due to the minimal computational cost of training simple pricing models from scratch, it is often **faster and more efficient** to train a new, dedicated model on the client's specific data rather than attempting complex fine-tuning.
- **Necessity of Fine-Tuning AI Agents:** For complex, real-world tasks, relying solely on **prompts over a foundational LLM** is unreliable and inefficient, especially when an agent must **decompose** a goal into multiple **actions**. Authentic AI agents require **fine-tuned**

**models** for most of their individual actions to ensure robust, reliable, and specialized behavior.

- **Fine-Tuning vs. In-Context Learning:** Fine-tuning wins over prompt engineering due to **reliability and deep character integration**. Prompting leads to **In-Context Learning**, where the model "reads the script" every time, resulting in **unreliable output** and a high error rate. Fine-tuning, through repeated training, bakes the desired behavior and specialized responsibility into the model's parameters, creating an **ingrained habit** with a **low error rate**.
- **Specialized Responsibility:** Fine-tuning creates a **specialized model** deeply imbued with a specific task (e.g., spotting regulatory concerns in a park plan). This deep integration is the solution to the biggest criticism against AI agents: their inconsistency, which stems from relying on prompts rather than specialized, fine-tuned models. Fine-tuning is a **last-mile problem**, making specialization relatively easy.
- **Cost-Benefit of Sizing the Brain:** Choosing the right model size involves a complexity study to find the **sweet spot** before performance suddenly **dips** when scaling down from a large **frontier model**. **Small models** (e.g., 2B) are used for simple tasks like junk filtering, while **Large Models** (e.g., 200B+) are necessary for complex tasks requiring **reasoning, planning, and task decomposition**.
- **Analogy of Specialization:** Using a fine-tuned model is like hiring a **specialist (neurosurgeon)** whose deep, ingrained expertise requires **minimal instructions (short prompt)**, leading to higher reliability. Conversely, a frontier model with a long prompt is like a **generalist** who needs a **mile-long set of instructions** every time.
- **The Escalation Ladder:** Application development should follow an **escalation ladder**, stopping wherever evaluation metrics are met: 1) System Prompt + LLM, 2) **Sidecar Mode (RAG Pipeline)** for fluid knowledge, 3) **Fine-tuning Models** for specialized reasoning/behavior change, and 4) Multi-Agent Reinforcement Learning (MARL).
- **RAG vs. Fine-Tuning is a False Dichotomy:** The debate is misleading; a good RAG system often *requires* fine-tuning of its components. **RAG** (Retrieval-Augmented Generation) is superior for **knowledge retrieval** (fluid knowledge base), while **Fine-tuning** is for changing the model's **behavior, reasoning, and specialized ability**.
- **Fine-Tuning is for Behavior, Not Knowledge:** It is mathematically illiterate to use fine-tuning to **bake knowledge** into a model's weights because that knowledge is stored in a **compressed, lossy form** with a high error rate upon retrieval. **RAG offers absolute preservation of knowledge** and is the reliable method for retrieval. Fine-tuning enhances the model's **conceptual understanding** so it can apply reasoning.
- **LLMs vs. Fine-Tuning Philosophies:** The debate is between the "**Larger is Better**" philosophy, which relies on **scale and emergent abilities** of large frontier models (LLMs), and the "**Specialization and Distribution**" strategy, which advocates for using a **web of smaller, fine-tuned models** (SLMs) within an agentic system.

- **Argument for Large LLMs:** LLMs benefit from **Neural Scaling Laws**, exhibiting **emergent abilities** and superior **domain mastery** simply due to their vast size (e.g., GPT-4 surpassing Bloomberg's specialized LLM). They offer architectural simplicity, handling complex tasks like reasoning and planning, and eliminate large upfront hardware costs (CAPEX) by using a pay-per-token structure (OPEX).
- **Argument for Fine-Tuned SLMs:** A **web of specialized, fine-tuned models** can collectively achieve or exceed the abilities of a single frontier model (**Distributed Intelligence**). They are far more economical for high-volume, repetitive tasks because, after initial CAPEX, the operational inference cost is near **zero** (no recurring token cost, OPEX advantage).
- **The Specialization Trap (MoE):** Relying on a giant model for a specialized task means you often pay to run infrastructure for the entire model (e.g., 600B parameters), when it may only activate one internal expert (e.g., 30B parameters). It is more efficient to **fine-tune a smaller 30B model in-house** to be the ultimate guru for that specific task.
- **Key Trade-Offs (Cost and Specialization):** The choice boils down to cost structure (LLMs use OPEX / token vs. Fine-Tuned use CAPEX / hardware), and specialization (LLMs are highly **adaptable/general**, while Fine-Tuned models are more **specialized/reliable** for a single task). Fine-tuned models also tend to run **faster** with lower latency.
- **A2A Enables Distributed Architecture** The A2A protocol is designed for a distributed architecture where agents run as independent servers, unlike single-process frameworks like CrewAI. This allows for true multi-agent communication across different companies and systems.
- **Clients Are Not Agents** In an A2A system, the application that initiates a request, such as the `A2A_SimpleClient`, is not an agent. It is simply a "client," analogous to a REST client, that sends messages to the actual agents performing the work.
- **Discovery is a Core Challenge** Agent discovery is a critical challenge, moving from basic Peer-to-Peer (P2P) where addresses are known in advance. The future vision is a centralized discovery model, like a "DNS for agents," to find unknown agents based on their capabilities.
- **The Future is Standardization** The current "protocol wars" between standards from Google, IBM, and Cisco are converging toward a single, open standard for agent communication. Google's A2A protocol, for example, has already been moved to the Linux Foundation to facilitate this convergence.
- **Move Beyond Simple ReAct Loops** To build reliable agents, advanced patterns like **Re-planning** and **Reflection** are necessary to move beyond simple ReAct loops. These patterns allow agents to recover from failure by creating new plans and to self-critique the quality of their own outputs.
- **The Goal is Reasoning, Not Just Logic** The "Smart Home" project illustrates the key value of AI agents: replacing cheap but rigid `if-then-else` logic with flexible, context-

aware reasoning. This allows the system to handle complex scenarios that traditional logic cannot, such as the "Sophie the Dog" analogy.

- **Industry Adoption is Still Early** The majority of companies (70%+) are not yet implementing these complex agentic patterns, as most enterprise work is still in the PoC phase. True, complex systems are currently limited to companies whose core business is AI, like Perplexity or Cursor.
- **The Industry is Shifting from "Generative" to "Reasoning" AI** The industry is evolving from just "Generative AI" to "Reasoning AI," which demands a new set of advanced skills. Key emerging skills include data preparation, advanced RAG, specialized model fine-tuning, and agentic orchestration.
- **The Fundamental Dichotomy:** Fine-tuning is for *Ability* (e.g., reasoning), while RAG is for *Knowledge* (e.g., facts). This separation is the most critical architectural principle.
- **LLMs are Not Knowledge Containers:** Trying to bake knowledge into an LLM is a flawed "ocean in a bucket" approach. Knowledge expands exponentially, while models are finite.
- **Fine-Tuning is Lossy Compression:** LLMs "grok" (generalize) by losing specifics when their memory capacity (3.6 bits/parameter) is full. This is unacceptable for recalling exact facts.
- **RAG Provides Provenance:** RAG systems can legally and verifiably cite their sources (provenance). LLMs relying on internal memory will *hallucinate* sources.
- **Fine-Tuning is Neurosurgery:** The process is high-risk, like neurosurgery; it's irreversible and can cause "catastrophic forgetting" or "mode collapse".
- **Good RAG Uses Fine-Tuned Models:** "All good RAG systems" use models fine-tuned for *ability* (like better reasoning), not for storing knowledge.
- **Local Agents Provide Governance:** A system of smaller, specialized agents provides traceability and governance. A single, giant black-box model removes this observability.
- **AI is About System Architecture:** Success is determined by the entire system architecture (embedders, rerankers, etc.), not one single model.
- **Use the Escalation Ladder:** Start with the simplest solution (e.g., system prompt) and only escalate to RAG or fine-tuning when necessary to meet business goals.
- **Apply Amdahl's Law:** Identify your system's weakest link (e.g., bad embeddings) and apply fine-tuning there for the biggest impact.
- **Frontier vs. Local (Cost):** Frontier models have zero **CapEx** (capital) but high, recurring **OpEx** (operational) costs. Local models are the reverse.
- **Interfaces Can Overcome Illiteracy:** Voice and camera interfaces can make advanced technology accessible to illiterate populations, removing illiteracy as a handicap.
- **Technology in Service of Mankind:** Engineering's goal, borrowing from the MIT motto, should be to create technology that provides tangible social utility and serves humanity.