

# Week2 ML Strategy

## Carrying Out Error Analysis

Look at dev examples to evaluate ideas



Should you try to make your cat classifier do better on dogs? ↩

- Error analysis: → 5-10 min
- { Get ~100 mislabeled dev set examples.
  - Count up how many are dogs.

$$\begin{array}{ll} \text{→ 5\%} & \text{10\%} \\ \underline{5/100} & \downarrow \quad \text{95\%} \end{array} \quad \left| \quad \begin{array}{l} \text{→ 50\%} \\ \text{50/100} \end{array} \right. \quad \begin{array}{l} 10\% \\ \downarrow \\ 5\% \end{array}$$

Andrew Ng

## What is Error Analysis?

- **Error analysis** is the process of manually examining the mistakes made by a machine learning model to gain insight into what improvements may be most effective.
- It's especially useful if your model is not yet performing at human-level accuracy on a given task.

## Why Do Error Analysis?

- Helps you decide where to focus your efforts for maximum gain.
- Rather than guessing or spending months on a potentially low-impact change, you can quickly estimate the benefit of fixing certain types of errors.

## Practical Example: Cat Classifier

- Suppose you have a cat classifier at **90% accuracy** (10% error) on your dev set.
- If you notice that your classifier often mistakes **dogs** as cats, should you focus on fixing that problem?
  - Instead of jumping in, **count** how many of the misclassified (error) examples are actually dog images.
- For example, if out of **100 error cases**, only **5 are dogs**:
  - **Fixing the dog problem** would at best improve your error rate from 10% to 9.5%.
  - **Result:** Not a huge gain; may not be worth major effort.
- If instead, **50 out of 100 are dogs**:
  - **Fixing** this could cut your error rate from 10% to 5%.
  - **Result:** A much more promising direction.

## Upper Bound (Ceiling) Concept

- The “**ceiling on performance**”: The maximum improvement possible by fixing a particular category of error.

## Error Analysis Procedure

- **Step 1:** Collect about 100 misclassified dev set examples.
- **Step 2:** Examine them manually and categorize each error (e.g., dogs, blurry images, great cats).
- **Step 3: Count** how many errors fall into each category.
- **Step 4:** The fraction gives you an estimate of the *potential improvement* if you fix that error type.
- **Step 5:** Use this analysis to decide if focusing on a category is worth significant effort.

## Analyzing Multiple Ideas in Parallel

# Evaluate multiple ideas in parallel

Ideas for cat detection:

- Fix pictures of dogs being recognized as cats ←
- Fix great cats (lions, panthers, etc..) being misrecognized ←
- Improve performance on blurry images ← ↘

Image	Dog	Great Cats	Blurry	Instagram	Comments
1	✓			✓	Pitbull
2			✓	✓	
3		✓	✓		Rainy day at zoo
⋮	⋮	⋮	⋮	⋮	
% of total	<u>8%</u>	<u>43%</u>	<u>61%</u>	<u>12%</u>	

Andrew Ng

- If you have several improvement ideas:
  - Make a **table or spreadsheet**, with:
    - Rows = Each misclassified example.
    - Columns = Each idea/category (dogs, great cats, blurry images, Instagram filters / snapchat filters,.. ,etc.)
    - Optional comments column for special notes.
- Go through each example and **check off** which errors apply.

## Adding New Categories

- While going through examples, you might discover new mistake categories (e.g., Instagram/Snapchat filters disturbing accuracy).
- It's fine to add new columns mid-way.

## Quantitative and Qualitative Insights

- At the end, sum up the percentages for each column (error category).
- Focus on categories with the highest impact potential.
- The process helps in deciding **how to prioritize** engineering efforts.

## Flexibility and Iteration

- Not a strict formula; use the findings to **guide decisions** rather than dictate them.
- Sometimes splitting work between teams for different high-potential improvements makes sense.

## Summary Process

1. Find a set of misclassified (error) examples in your dev set.
  2. For each, categorize into possible mistake types.
  3. Count the number/fraction in each category.
  4. Estimate which fixes will give the most performance gain.
  5. Use as a **fast, evidence-based method** to guide future work instead of blind guessing.
- 

### Key Takeaways:

- **Manual error analysis** lets you estimate the maximum gain possible by fixing various types of errors.
- **This helps prioritize your efforts** for maximum impact, saving time and resources.
- The process is **fast** (5-10 minutes for 100 errors) and **flexible**—add new categories as needed.
- **Apply error analysis regularly** as your model and problem domain evolve.

## Cleaning Up Incorrectly Labelled Data(in the dataset )

# Incorrectly labeled examples



DL algorithms are quite robust to random errors in the training set.

## Error analysis

Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
98				✓	Labeled missed cat in background
99		✓			
100				✓	Drawing of a cat; Not a real cat.
% of total	8%	43%	61%	6%	

Overall dev set error ..... 100%

Errors due to incorrect labels ..... 0.6%

Errors due to other causes ..... 9.4%

2%  
0.6%  
1.4%  
2.1%  
1.9%

Goal of dev set is to help you select between two classifiers A & B.

## 1. Definition of Incorrectly Labeled Data

- **Incorrectly labeled data:** The ground-truth label ( $Y$ ) assigned by the human labeler is wrong (e.g., a dog labeled as "cat").
- Distinction:
  - *Mislabeled example:* Model predicts wrong label ( $Y$ ).
  - *Incorrectly labeled example:* Actual dataset label is wrong.

## 2. Handling Incorrect Labels in the Training Set

- Deep learning models are robust to random errors in the training set labels (e.g., labeler pressed the wrong key by accident).
  - If errors are random and relatively rare, it's often okay to leave them—model performance is usually not seriously harmed.
  - No harm in fixing labels if time permits, but not always necessary.
  - Systematic errors (e.g., always labeling white dogs as cats) are more harmful than random ones. These can cause the model to learn incorrect rules—must be fixed.
  - As long as dataset is large and errors are few and random, model performs well.
- 

## 3. Incorrect Labels in Dev/Test Sets—Why It Matters

- Dev/test labels must be accurate—these sets are used for:
    - Model selection
    - Hyperparameter tuning
    - Final evaluation
  - Incorrect labels cause incorrect error measurements, misleading your assessment.
  - During error analysis, keep track of misclassified examples caused by wrong labels.
    - Add an extra column in your error analysis for tracking incorrect-labeled dev/test set samples.
- 

## 4. Prioritizing Label Correction in Dev/Test Sets

- Assess the impact of incorrect labels:
  1. Dev set error rate (e.g., overall 10% error).
  2. Fraction of errors due to wrong labels (e.g., 6% of errors due to mislabeled data).
  3. Errors due to all other causes.

- **Example:** If 0.6% of 10% error is due to label mistakes, you have 0.6% error from labeling, 9.4% from other issues.
  - If overall error falls (e.g., to 2%), but label mistakes remain 0.6%, now 30% of errors come from label errors—makes label fixing more worthwhile.
- 

## 5. When to Fix Incorrect Labels

- Fix dev/test labels if:
    - They significantly impact your ability to select/tune models.
    - The fraction of error due to bad labels is high.
  - If fixing will not change selection/evaluation, it may not be best use of time.
- 

## 6. How to Fix Labels—Best Practices

- **Apply label correction process to both dev and test sets at once** to keep distributions identical.
  - Examine *both*:
    - Cases the model got wrong (false positives/negatives).
    - Cases the model got right (could be right by chance when label is wrong).
  - It's harder (time-consuming) to check "right" cases—can prioritize "wrong" ones with tradeoffs.
- 

## 7. Training Set Label Correction

- **Less important** to fix small numbers of incorrect labels in a large training set.
  - Dev/test sets, being smaller, are more manageable for manual relabeling.
- 

## 8. Real-World Practice and Advice

- Modern deep learning often favors large-scale, less manual processing.
- However, **manual error analysis and some human insight remain important**, especially for understanding and prioritizing improvements.

- Even accomplished researchers manually count and analyze errors for critical projects—*this is valuable time spent.*
- 

## 9. Summary Takeaways

- **Random label errors** in big training sets are tolerable.
  - **Systematic label errors** (patterns) **must be fixed**.
  - **Dev/test labels must be fixed** if they distort model evaluation.
  - **Manual error analysis** (including label review) is a key part of successful machine learning projects—even in the deep learning era.
- 

# Build your First System Quickly, then Iterate

## Speech recognition example



- • Noisy background
    - • Café noise
    - • Car noise
  - • Accent
    - • Far from native
    - • Young
    - • Stutter
    - • ...
- Guideline:  
**Build your first system quickly, then iterate**
- • Set up dev/test set and metric
  - Build initial system quickly
  - Use Bias/Variance analysis & Error analysis to prioritize next steps.

## 1. Philosophy: Embrace Iteration, Not Perfectionism

- Building a new ML application often has *many possible improvement directions* (e.g., robustness to noise, handling accents, special cases like children's speech, etc.).
- Even experts can't always know ahead of time which improvement will matter most in a new domain.

- The key question: “**What should I prioritize?**”
- 

## 2. Step 1: Build a Quick, Working Version

- **Rapid prototyping:** Quickly set up your dev/test set and choose an appropriate metric. Don’t worry about finding the perfect metric now—you can adjust it later.
  - Pull together an initial version of your system—“quick and dirty” is acceptable.
  - Main goal: **Get any system working** so you can start measuring its performance.
- 

## 3. Step 2: Use Analysis to Target Improvements

- Once your first attempt is running:
    - Apply **bias/variance analysis** to see if you’re underfitting or overfitting.
    - Perform **error analysis**: Review mistakes, categorize error types, and identify which issue is most critical.
  - Example: If most errors are from far-field (distant microphone) voices, target “far-field speech” techniques.
  - This analysis **localizes the problem** and helps you invest effort where it will give the most return.
- 

## 4. Step 3: Iterate

- Now, *target the biggest issues* you’ve found with concrete improvements.
  - Repeat the loop:
    - Implement a fix.
    - Analyze results.
    - Prioritize next experiment based on data.
  - Continue until performance is satisfactory.
- 

## 5. Why This Approach Works (Empirical Mindset)

- Scientific mentality: Let **data, not speculation, guide your next step.**
  - Avoids “analysis paralysis” from over-planning or guesswork—iteration reveals which challenges matter most.
  - Real-world machine learning almost always requires adaptation on the fly.
- 

## 6. Exceptions to the “Quick First System” Rule

- If the application has a deep academic background and your problem matches it closely (e.g., face recognition), you can build on existing, complex systems with more confidence.
  - If you already have **significant prior experience** in the domain, you may be able to skip some exploratory steps.
- 

## 7. Common Pitfalls & Advice

- **Overthinking:** Many teams build unnecessarily complex systems before any data analysis.
  - **Underthinking:** Some teams build overly simple systems but this is less common—mistakes can be corrected via iteration.
  - **Best practice:** On average, it's safer to start simple and iterate, than start complex without empirical justification.
- 

## 8. Bottom Line / Practical Advice

- For innovation-focused applied ML: **“Build something that works—quickly —then iterate and learn from your data.”**
  - Each prototype and experiment informs your next move; stay agile and objective.
  - Error and bias/variance analysis are your scientific tools for progress.
- 

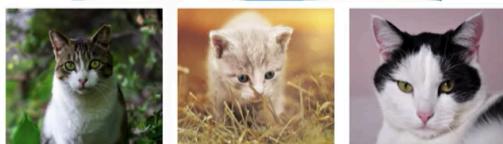


Build your first system quickly, then iterate

# Training and Testing on Different Distributions

Cat app example

Data from webpages



$\rightarrow \approx 200,000$

$210,000$   
shuffle

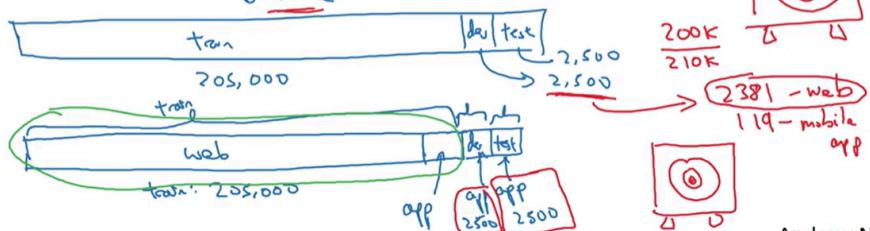
core about this

Data from mobile app



$\approx 10,000$

X Option 1:



Option 2:

app  
2,500

Andrew Ng

Speech recognition example

Speech activated rearview mirror



Training

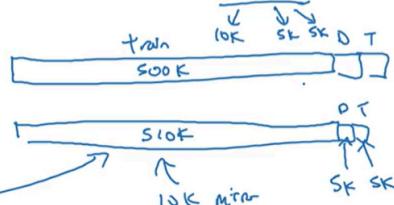
- { Purchased data  $x, y$
- Smart speaker control
- Voice keyboard

...  
 $500,000$  utterances

Dev/test

Speech activated  
rearview mirror }

$\rightarrow 20,000$



## 1. Core Concept:

When building machine learning systems, the data you use for training may often come from a different distribution than the data your model sees in

practice (dev/test time). The key idea is to set up your data splits so your system learns from all available data—even if some is “different”—but evaluates solely on the real-world data you care about, so your results are meaningful and actionable.

---

## 2. Breakdown:

### How It Works:

- Gather all available data: Often, you may have lots of “easy to get” data (e.g., web images, industry datasets) and a smaller set that actually matches your deployment setting (e.g., user-uploaded mobile photos).
- Split wisely:
  - **Training set:** Use both “source” (different distribution) and “target” (real-world) data to maximize learning.
  - **Dev/Test sets:** Use only the real-world (target) data, so your model is tuned and assessed for the conditions it will face in actual use.
- Example Use Cases:
  - Cat detector: Training on web + app images, but testing only on app uploads.
  - Speech system: Train on various speech corpora + new device utterances, but evaluate only on new device queries.

### Key Formulas/Quantitative Intuition:

- If you pool all data and split randomly, the dev/test set may disproportionately represent the “easy” (non-target) source, invalidating your performance estimates.
- Instead:
  - Train set = Source data + Part of Target data
  - Dev/Test sets = Exclusively Target data
- For example, if you have 200,000 source + 10,000 target images, use all 210,000 for training, but put all 2,500 dev/test samples from the target.
- Variables:
  - $n_{\text{source}}$  = source samples (web, generic)
  - $n_{\text{target}}$  = target samples (mobile app, real-world)

- $n_{\text{train}}, n_{\text{dev}}, n_{\text{test}}$  = respective splits

### Why It's Used:

- **Problem solved:** Makes the most of all available data while ensuring the model is tuned for real-world use.
  - **Advantages:**
    - Avoids setting false goals by overfitting to easily available but irrelevant data.
    - Yields more honest model evaluations and real-world deployment readiness, especially when "distribution mismatch" arises.
- 

### 3. Practical Insights:

#### Real-World Applications:

- **Face Recognition:** Training on diverse celebrity datasets, evaluating on security camera images from your building.
- **Speech-activated Devices:** Training on big open speech datasets + limited data from specific accents/devices; testing on actual user utterances.
- **Medical Diagnosis:** Model trained on open-source scans from various hospitals, but validated and tested only on local patients' scans.

#### Common Pitfalls:

- **Dev/test contamination:** Don't let "source" data creep into dev/test—your numbers will be misleading!
  - **Overfitting to irrelevant features:** Mixing non-target data in dev/test causes the model to optimize for the wrong patterns.
  - **Not enough target data in train:** Too little relevant data in train means the model won't generalize to real usage. Consider data augmentation if possible.
- 

### 4. Creative Extension:

#### Project Idea:

Build an object recognition system for food items in canteen photos. Collect open datasets (e.g., stock food images) and combine with a small set of photos snapped in your local canteen on your phone. Train on both, but use only those

canteen photos for dev/test. Analyze how performance differs on “clean” vs. “real” images and try techniques to close the gap.

#### Open Question:

How can we best select or “weigh” source samples so they’re maximally helpful for learning, without introducing bias toward the source distribution? Could a learned weighting or domain adaptation help bridge the gap more systematically?

---

**Summary tip:** Always ask: “What data does my system really need to succeed on in production?” Your dev/test splits should reflect this—even if your train data is a wild mix!

## Bias and Variance with Mismatched Data Distribution

### Data Mismatch:



**Data mismatch** refers to a situation where the data used to **train your model** is different in distribution or characteristics from the data you actually encounter in production or during evaluation (dev/test sets).

#### In plain terms:

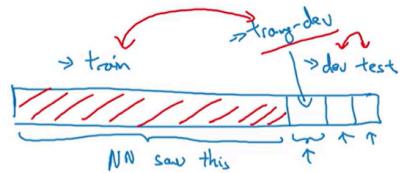
- You train your AI on one kind of data.
- But when you test it (or use it in real-world situations), the data looks or behaves differently.
- Because of this difference, your model’s performance on the training data does **not accurately predict** its performance in the real-world scenario.

## Cat classifier example

Assume humans get  $\approx 0\%$  error.

Training error ..... 1%  $\downarrow$  9%  
 Dev error ..... 10%

Training-dev set: Same distribution as training set, but not used for training



	Training error	1%	Variance	1%
$\rightarrow$ Training-dev error	9%		1.5%	
$\rightarrow$ Dev error	10%		10%	
				Variance
Human error	0%	↑ Available bias	10%	↑ Available bias
Training error	10%		11%	
Training-dev error	11%		11%	
Dev error	12%		20%	
		Bias		↑ Variance
				↑ Data mismatch
				↑ Bias + Data mismatch

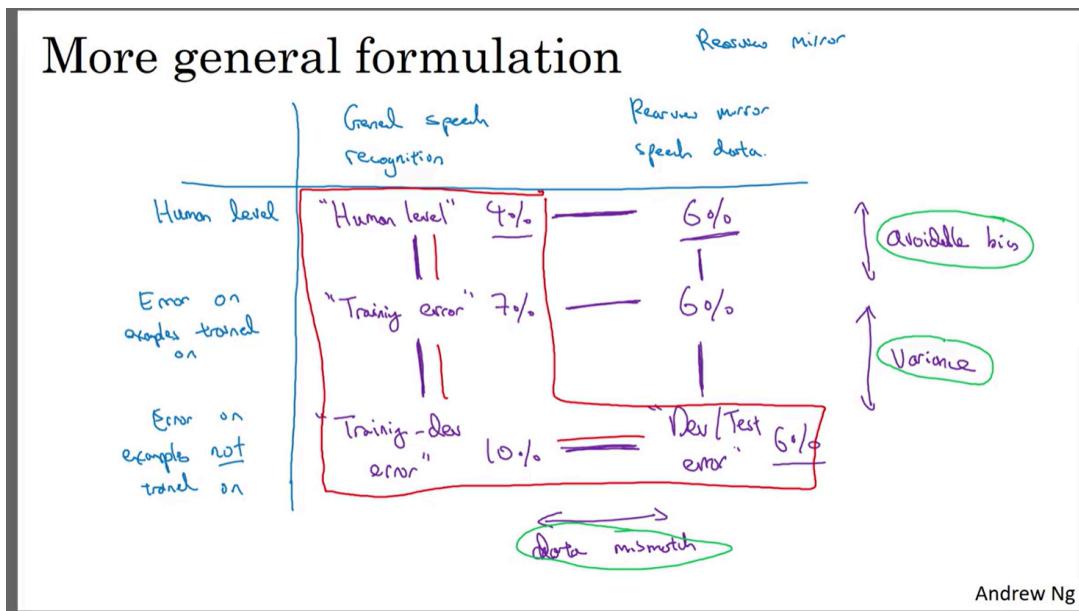
Andrew Ng

## Bias/variance on mismatched training and dev/test sets

Human level	4%	7%	10%	12%	12%
Training set error					
Training-dev set error	4%	7%	10%	12%	12%
$\rightarrow$ Dev error					
$\rightarrow$ Test error					
	↑ available bias	↑ variance	↑ data mismatch	↑ degree of similarity to dev set	

4%
7%
10%
6%
6%

## More general formulation



### 1. Core Concept

When your training data comes from a different distribution than your dev/test data, it's not enough to just look at training and dev error: you need more nuanced analysis (introducing "training-dev" set) to properly disentangle bias, variance, and data mismatch problems.

### 2. Breakdown

#### How It Works:

- Normally, we analyze bias and variance by comparing performance on the training and development (dev) sets, assuming they share the same distribution.
- If the distributions differ, a direct comparison can mislead: a big dev set error might be due to distribution differences, not just lack of generalization.
- Solution:
  - **Training Set:** Data used to fit the model.
  - **Training-Dev Set:** Subset from the training set distribution not seen during training (never used to update model weights).
  - **Dev Set:** From the target (production) distribution.
  - Compute errors on all three.
- Analysis:

- **Gap between Training and Training-Dev error:** Measures variance (generalization within the training distribution).
- **Gap between Training-Dev and Dev error:** Measures data mismatch (distribution shift).
- **Training error vs Human-level error:** Measures avoidable bias.

### **Key Formulas:**

- Let,
  - $E_{\text{human}}$  : Human-level (Bayes error, ideal performance)
  - $E_{\text{train}}$  : Error on training set
  - $E_{\text{train-dev}}$  : Error on training-dev set
  - $E_{\text{dev}}$  : Error on dev set

$$\text{Bias} \approx E_{\text{train}} - E_{\text{human}}$$

$$\text{Variance} \approx E_{\text{train-dev}} - E_{\text{train}}$$

$$\text{Data Mismatch} \approx E_{\text{dev}} - E_{\text{train-dev}}$$

- Each variable is an error rate; the differences quantify the source of poor generalization.

### **Why It's Used:**

- **Problem solved:** Figure out why your model underperforms—whether it's due to underfitting, overfitting, or the data being different in training and production.
- **Advantages:**
  - *Pinpoints root causes more accurately than basic bias-variance split.*
  - *Helps you prioritize whether to get more/better data, regularize, or improve modeling.*

## **3. Practical Insights**

### **Real-World Applications:**

- *Autonomous driving:* Training on clean/ideal images, but dev/test from real road conditions (night, rain, dirt).

- *Voice assistants (ASR)*: Training on lab speech data, but real users/dialects in deployment.
- *Medical imaging AI*: Training data from one hospital/device, deploying to others with different equipment or demographics.

### **Common Pitfalls:**

- *Ignoring Data Mismatch*: Only comparing training and dev error can fool you—your changes (regularization, more training, etc.) might not address the real issue.
- *Underestimating the need for training-dev set*: Without this, you might misdiagnose variance when the real problem is distribution shift.
- *Not measuring human-level performance*: Without this, you can't judge if error is even avoidable.

### *How to Avoid Them:*

- Always create a training-dev set and evaluate on it.
- Track and write down human-level or Bayes-optimal error when possible.
- When large gaps are found:
  - *Bias problem*: Improve model/feature engineering.
  - *Variance problem*: Get more data/regularize.
  - *Data mismatch*: Collect data more like your target or use domain adaptation methods.

## **4. Creative Extension**

### **Project Idea:**

Build a mini image classifier (e.g., for cats vs. dogs) where:

- Training set: High-res, “studio” photos.
- Training-dev set: Unseen studio images.
- Dev set: “In the wild” phone captures.

Run the full error analysis and identify which problem is largest (bias, variance, mismatch). Try to improve performance on the dev set by applying suitable strategies—document your findings.

### **Open Question:**

What are systematic, general-purpose methods for solving data mismatch (domain shift) in deep learning that work reliably across different domains and data types? Many such methods exist, but none are broadly robust—

## Addressing Data Mismatch

Example Car RearView mirror

- • Carry out manual error analysis to try to understand difference between training and dev/test sets

E.g. noisy - car noise      street numbers

- • Make training data more similar; or collect more data similar to dev/test sets

E.g. Simulate noisy in-car data

## Artificial data synthesis



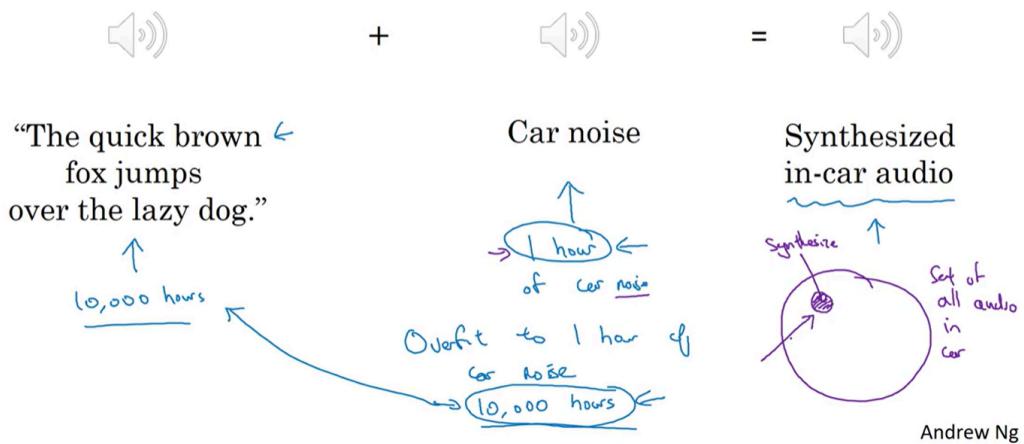
“The quick brown fox jumps over the lazy dog.”

Car noise

Synthesized  
in-car audio

mixing with what data you have with car noise is better idea for synthesis

# Artificial data synthesis



might overfit if the noise is similar and look like as small subset

Car recognition:



for car recognition we can ai data synthesis but not always good idea it might overfit

## 1. Core Concept:

Data mismatch happens when your **training set** is drawn from a different distribution than your dev/test set. This difference can dramatically hurt model performance, especially when real-world production data doesn't match your clean training data. Addressing this mismatch is crucial to create robust, high-performing AI systems.

## 2. Breakdown:

### How It Works:

- **Identify the Problem:** Use error analysis to spot if your model struggles because of a mismatch. Focus on dev set errors, not test set, to avoid overfitting.
- **Understand the Nature of Mismatch:** Look for patterns (e.g., extra noise, new categories, unfamiliar accents in speech tasks) present in dev/test but not in training.
- **Make Training Data More Realistic:**
  - **Artificial Data Synthesis:** Manually or programmatically inject features (noise, special cases) of dev/test into training samples.
  - **Targeted Data Collection:** Gather more samples resembling the challenging cases in your dev/test set.

### **Algorithm Steps / Key Methodology:**

1. Perform manual error analysis on the dev set to discover systematic mismatches.
2. Simulate or collect data resembling the dev/test distribution.
  - Example: For speech recognition in cars, mix clean training voice data with real car noise backgrounds.
3. Cautiously expand your dataset, ensuring generated data covers a wide range, not just a narrow slice.

### **Key Formulas:**

There are no complex formulas—this process is more about data engineering than mathematics. However, if synthesizing data, the basic formula for audio mixing:

$$\text{synthetic\_audio}(t) = \text{clean\_audio}(t) + \text{noise}(t)$$

Where:

- $\text{clean\_audio}(t)$  : original audio at time  $t$
- $\text{noise}(t)$  : background noise signal at time  $t$

### **Variable Intuition:**

- **clean\_audio:** The original, noise-free sample (what your model has already seen).

- **noise:** The real-world effect you want to mimic (e.g., car sounds, reverberation).

### Why It's Used:

- **Problem Solved:** Bridges the gap between overly clean/ideal training data and messy, realistic data encountered in actual deployment.
  - **Advantages:** Dramatically improves performance on real-world data, prevents overfitting to irrelevant artifacts, and helps anticipate edge cases.
- 

### 3. Practical Insights:

#### Real-World Applications:

- *Speech Recognition in Cars:* Training with artificially noised audio to ensure the system works even with car/road sounds.
- *Self-driving Cars:* Using synthetic computer graphics to generate images of cars in varied environments for robust detection.
- *Medical Imaging:* Simulating different scanner settings or patient backgrounds to ensure models generalize to new hospitals.

#### Common Pitfalls:

- **Overfitting to a Narrow Range:** If you only use a small set of background noises (e.g., a single hour of car noise), your model may overfit and perform poorly when facing unseen varieties.
- **Simulated Data Not Matching Reality:** Synthesized data that *looks* or *sounds* realistic to a human may still miss critical diversity, leading to hidden blind spots.
- **Test Data Peak:** Accidentally tuning approaches based on test set feedback can lead to over-optimistic benchmarks. Only use the dev set for error analysis.

#### How to Avoid:

- Use as diverse a set of synthesized and collected data as possible.
  - Alternate between different synthesis techniques and collect real-world samples if feasible.
  - Keep the test set untouched for a true measure of generalization.
- 

### 4. Creative Extension:

### **Project Idea:**

*Build a robust speech recognizer for “noisy environments”.*

- Collect/curate clean voice samples.
- Collect several distinct types of real-world noises (cafes, traffic, office, etc.).
- Synthesize training data by mixing clean and noise.
- Analyze performance improvements and discuss the impact of varying noise types.

### **Open Question:**

*How can we systematically measure or guarantee that synthesized data covers the “full diversity” of real-world scenarios, not just a visually/audibly convincing—but incomplete—subset?*

This remains an open challenge in ensuring data synthesis truly generalizes model behavior.

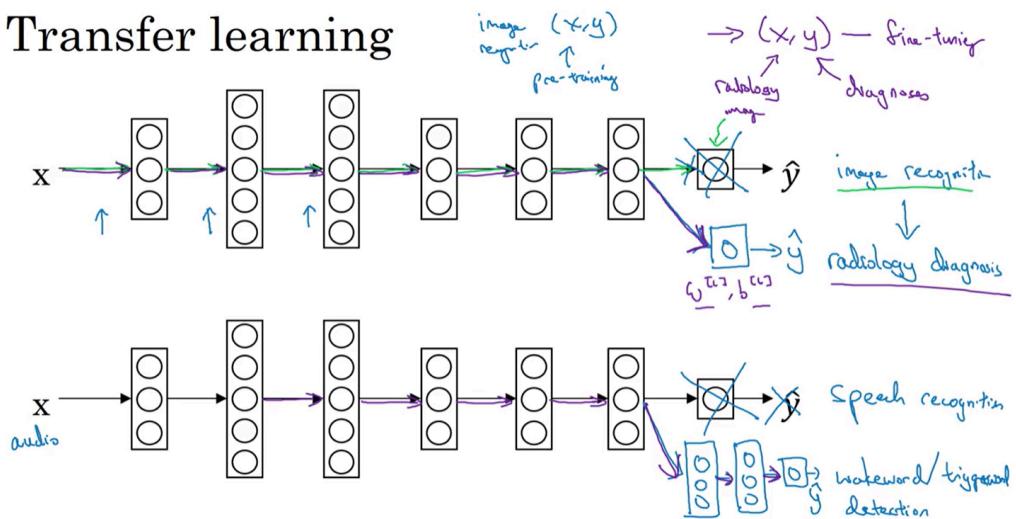
---

### **Summary:**

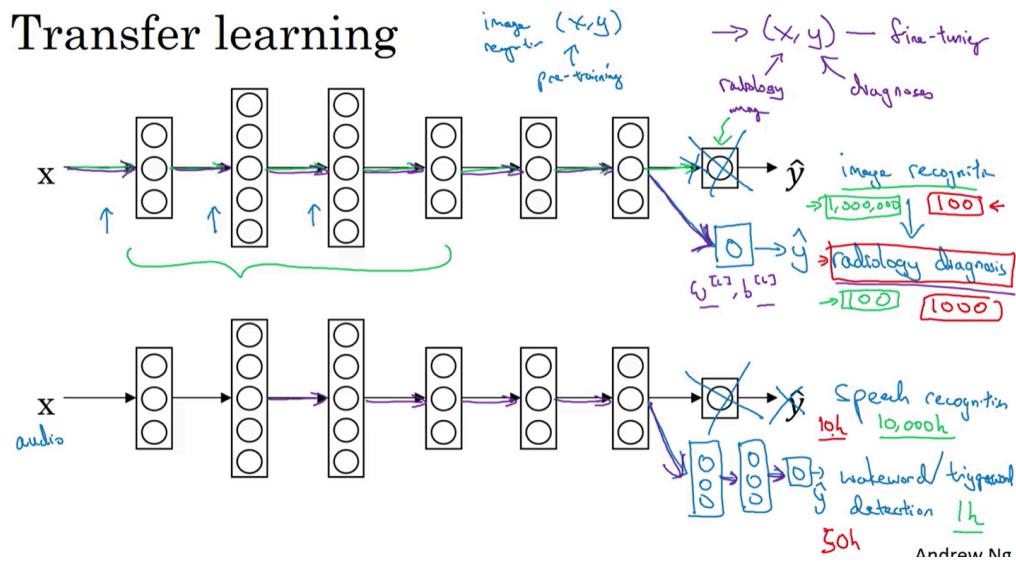
Addressing data mismatch involves recognizing, analyzing, and bridging gaps between synthetic/ideal and real-world data. While techniques like artificial data synthesis can make your models more reliable, they must be applied thoughtfully to avoid hidden overfitting and ensure your model is ready for the complexity of reality

## **Transfer Learning**

## Transfer learning



## Transfer learning



### 1. Core Concept

**Transfer learning** is a technique where you take the knowledge a neural network has learned from one task (Task A) and apply it to a different but related task (Task B). Instead of training a new model from scratch for every new task, you reuse useful features or "knowledge" learned previously, which accelerates training and improves performance—especially when you have less data for Task B.

### 2. Breakdown

#### How It Works

- **Step 1:** Train a neural network on Task A, typically with a large dataset (e.g., object recognition with millions of labeled images).
- **Step 2:** Remove the final layer(s) of the trained network (since the output classes usually differ between tasks).
- **Step 3:** Add a new layer(s) specific to Task B (e.g., radiology diagnosis or wake word detection).
- **Step 4:** Depending on the available data for Task B:
  - If **small dataset**: Retrain only the new final layer(s), freezing the rest.
  - If **large dataset**: Retrain (fine-tune) the entire network for the new task.
- **Step 5:** The earlier parts of the neural network retain learned "low-level" features (edges, curves for images, basic phonemes for audio) that are often transferrable between tasks with similar input structures.

## Key Formulas

- Let  $X$  = input (e.g., image or audio).
- Let  $f(X; W_{1..L-1})$  = output through the pre-trained layers.
- For Task A:  $Y_A = g_A(f(X; W_{1..L-1}); W_L^A)$
- For Task B: Replace output layer with new weights:
  - $Y_B = g_B(f(X; W_{1..L-1}); W_L^B)$
- Variables:
  - $W_{1..L-1}$  = shared, learned parameters.
  - $W_L^A, W_L^B$  = task-specific parameters for the final layer.

## Why It's Used

- **Solves the problem:** When you don't have enough data to train a deep model from scratch on a new task (Task B)—but similar data exists for a related problem (Task A) with ample examples.
- **Advantage:** The model leverages learned, general-purpose features from Task A, speeding up learning and boosting accuracy for Task B (where data is scarce). It prevents overfitting and reduces computational resources compared to training everything anew.

---

## 3. Practical Insights

## Real-World Applications

- **Medical Imaging:** Pre-trained image networks on generic visual data, then fine-tuned to classify diseases in X-rays or MRIs, where medical data is limited.
- **Speech Recognition:** Large speech-to-text models pre-trained on hours of general audio then adapted for limited-shot "wake word"/trigger word detection (e.g., "Okay Google," "Alexa").
- **Natural Language Processing (NLP):** Large language models like BERT/GPT trained on vast corpora, then specialized for sentiment analysis, translation, or legal document classification with much smaller labeled datasets.

## Common Pitfalls

- **Negative Transfer:** If Task A and Task B are not related enough, pre-learned features might hurt performance.
  - *How to avoid:* Only transfer when tasks share similar input structure/data.
- **Overfitting (when fine-tuning):** Especially if Task B's dataset is small.
  - *How to avoid:* Freeze most layers, retrain only the final layer(s); employ regularization, data augmentation.
- **Resource Wastage:** If Task B already has plenty of data or is drastically different from Task A, transfer learning may not be needed.

---

## 4. Creative Extension

### Project Idea

- *Build an image classifier for medical X-ray scans* from a small dataset by fine-tuning a classic image recognition model like ResNet (pre-trained on ImageNet). Evaluate how performance changes by unfreezing different layers and using various data augmentation techniques.

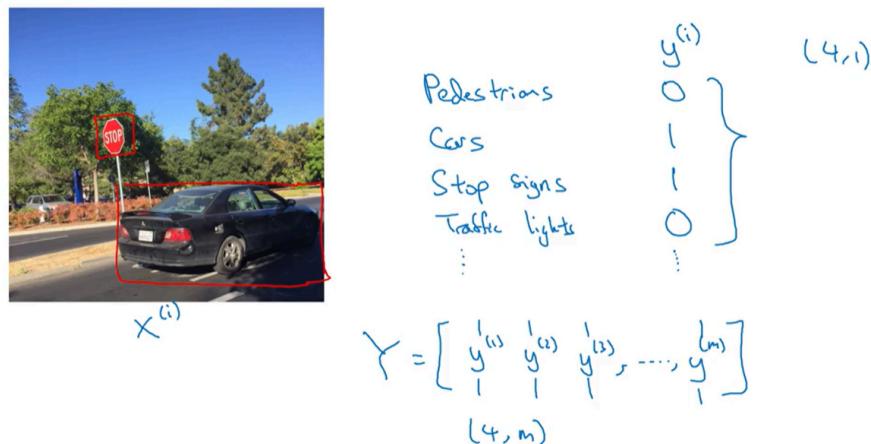
### Open Question

- *How can we better quantify and predict when transfer learning will lead to positive (vs. negative) transfer, especially as models and tasks get more complex?* An automated, reliable metric for "transferability" could save enormous time and resources in practice.

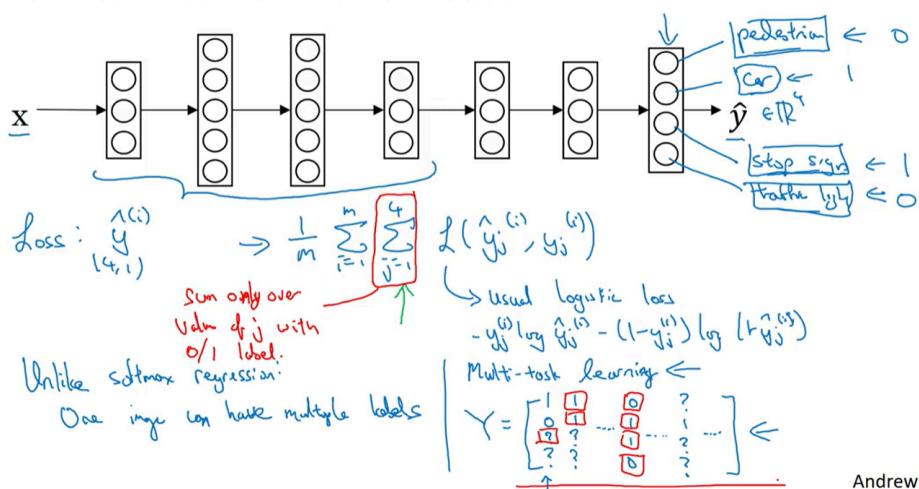
**Summary Tone:** As a peer, transfer learning is like "recycling" model knowledge for new, related problems—saving you time, data, and computation. It works best when you're moving from a "data-rich" area to a "data-poor" one, provided both problems share similar input types and data characteristics. Always check the relevance of tasks to avoid negative transfer, and use regularization methods to prevent overfitting on the scarce data for the new task.

## Multi-task Learning

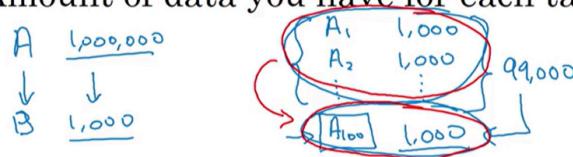
Simplified autonomous driving example



Neural network architecture



## When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.  

- Can train a big enough neural network to do well on all the tasks.

## Detailed Notes: Multi-task Learning (Deep Learning Specialization)

### 1. Core Concept

Multi-task learning (MTL) is a technique where a single neural network is trained to perform several related tasks simultaneously, leveraging shared features so that learning on one task can help the others. Instead of solving each task with separate models, MTL trains one model to solve them all together.

### 2. Breakdown

#### How It Works:

- In contrast to transfer learning (where knowledge is moved from one task to another, usually sequentially), MTL works by tackling multiple tasks in parallel.
- A neural network receives an input (e.g., an image) and predicts several outputs at once (e.g., for a self-driving car: pedestrian, car, stop sign, traffic light detection).
- Each output node (or group of nodes) corresponds to a specific task. The network's early layers learn shared features useful for all tasks, while later layers may specialize for individual task outputs.
- When labels are available for some but not all tasks (i.e., missing labels), the loss is simply computed for existing labels per sample, omitting missing ones from loss calculation.

## Key Formulas:

- For  $m$  examples and  $k$  tasks (labels), the loss can be written as:

$$\text{Loss} = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^k \mathcal{L}(y_j^{(i)}, \hat{y}_j^{(i)})$$

where:

- $y_j^{(i)}$  : actual label for task  $j$  on sample  $i$  (usually 0/1 for presence/absence)
- $\hat{y}_j^{(i)}$  : predicted probability for task  $j$  on sample  $i$
- $\mathcal{L}$ : loss per task, commonly logistic loss (for binary tasks):  
$$\mathcal{L}(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

## Why It's Used:

- **Problem:** In scenarios with related prediction tasks, doing them separately ignores shared information.
- **Advantages:**
  - Helps when tasks can benefit from shared features (e.g., images of roads for various traffic objects).
  - Allows smaller datasets for an individual task to be "amplified" by data from other related tasks.
  - Can improve performance and generalize better than training many isolated models, especially when data is balanced across tasks and the model is large enough.

---

## 3. Practical Insights

### Real-World Applications:

1. **Autonomous Vehicles:** Detecting multiple object types (cars, pedestrians, signs, traffic lights) in one forward pass through a single network.
2. **Medical Diagnosis:** Predicting multiple related conditions or risk factors from a common set of medical images or patient data.
3. **Natural Language Processing (NLP):** Jointly predicting sentiment, topic, and named entities in text using one network.

### Common Pitfalls:

- **Imbalanced data per task:** If some tasks have much more data, learning may be biased toward those.
  - **Model too small:** If the neural network doesn't have enough capacity, tasks can interfere (negative transfer), hurting overall performance.
  - **Unrelated tasks:** Forcing unrelated tasks into a shared model can degrade results (information sharing is only helpful if there is real feature overlap).
  - **Incomplete labels:** Datasets may have missing labels for some tasks; these should be omitted from the loss for those samples (handled by summing over only available labels).
- 

#### 4. Creative Extension

##### **Project Idea:**

Build a multi-label image classifier using an open dataset (e.g., multispecies plant/animal images) to predict several attributes at once (like species, color, presence of flowers/fruits). Experiment and compare against separate models for each attribute to observe the benefit of shared representations.

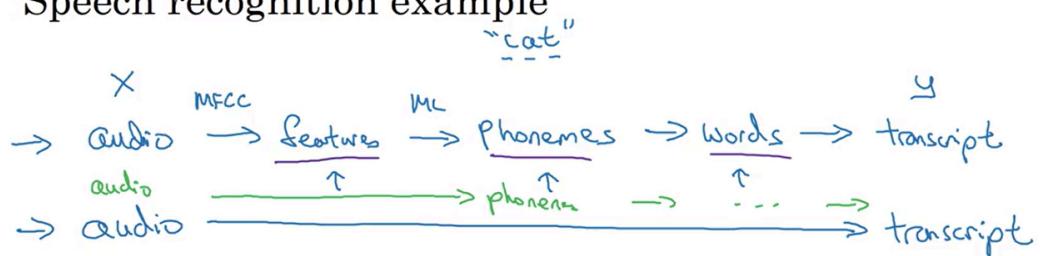
##### **Open Question:**

How can we design adaptive architectures or loss functions in multi-task learning that automatically detect and avoid negative transfer between tasks, especially for large-scale multi-task setups with partially unrelated targets?

---

## What is End-to-end Deep Learning?

## Speech recognition example



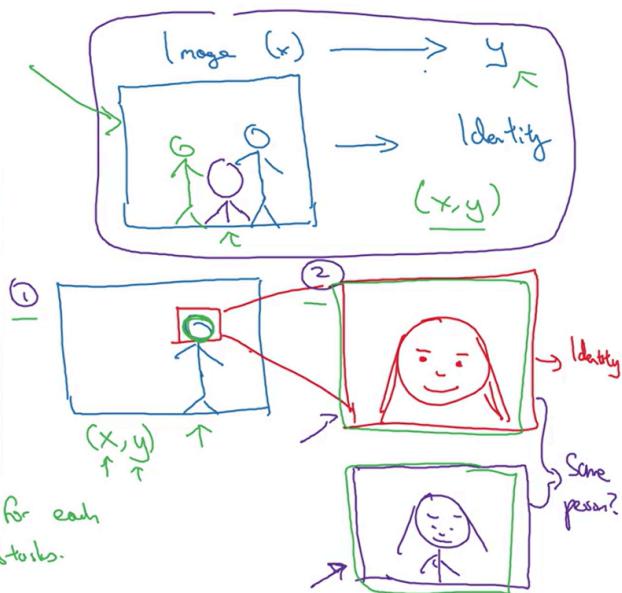
$3,000h$   
↑  
 $10,000h$   
⋮  
 $100,000h$

## Face recognition



[Image courtesy of Baidu]

Have data for each  
of 2 subtasks.



## Machine translation



Estimating child's age:



Image  $\xrightarrow{\textcircled{1}}$  bones  $\xrightarrow{\textcircled{2}}$  age

Image  $\longrightarrow$  age

Andrew Ng

## 1. Core Concept (Plain English Overview)

**End-to-end deep learning** is an approach where a single neural network is trained to map raw input data directly to the desired output, replacing traditional systems that use multiple sequential stages of manual feature extraction and processing.

## 2. Breakdown

### How It Works:

- Traditional systems involve several pipeline stages (e.g., for speech recognition: feature extraction  $\rightarrow$  phoneme detection  $\rightarrow$  word assembly  $\rightarrow$  transcription).
- End-to-end learning skips intermediate manual stages and trains one large model (usually a deep neural network) to learn the full mapping (input XXX  $\rightarrow$  output YYY).
- This approach can simplify pipelines, but often requires *large datasets* to outperform classical methods.
- Hybrid approaches exist, where only some stages are replaced by learned models.

### Key Formulas:

- Core mapping:  $y = f_\theta(x)$   
◦ xxx: raw input (audio clip, image, etc.)  
◦ yyy: desired output (transcript, label, etc.)

- $f_{\theta}$ : neural network with learnable parameters  $\theta$
- Loss function:  $L(f_{\theta}(x), y)$ , minimized during training to adjust  $\theta$

### **Variable Explanations:**

- $x$ : Represents the raw form of input (unprocessed data, like a sound waveform or an image).
- $y$ : The output the model should produce (text, category, number, etc.).
- $\theta$ : All the network's weights that are adjusted during training.
- $L$ : Quantifies the error between the model's guess and the true answer—the “training signal”.

### **Why It's Used:**

- **Problem Solved:** Eliminates dependency on hand-crafted features or intermediate modules, automates learning process end to end.
  - **Advantages:**
    - Simplifies the architecture—less manual engineering.
    - If enough data exists, can outperform systems with many manual steps.
    - Adapts flexibly to many different input/output relationships.
- 

## **3. Practical Insights**

### **Real-World Applications:**

- **Speech Recognition:** Raw audio is fed directly into a neural net, which outputs the transcript—skipping manual phoneme and feature extraction.
- **Machine Translation:** Networks learn to map whole sentences from one language to another directly, outperforming rule-based or multi-stage systems if large X-Y (sentence pair) datasets exist.
- **Medical Imaging:** For tasks like estimating a child's age from an X-ray, an end-to-end model could, in theory, go straight from the X-ray image to the age prediction—but data size limitations often still force multi-step systems.

### **Common Pitfalls:**

- **Data Hunger:** End-to-end deep learning *needs vast amounts of data*; with small or medium data, classical staged systems often outperform.

- **Uninterpretable Failures:** Bugs or issues are harder to track down since everything is crammed into one model.
  - **Overfitting:** With insufficient data, models may learn spurious correlations.
  - **Complexity vs. Data:** Sometimes, breaking tasks into sub-problems gives better accuracy due to available labeled data (e.g., face detection + recognition works better than straight image-to-identity if you have more labeled pairs for subtasks).
- 

## 4. Creative Extension

### Project Idea:

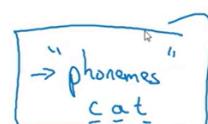
- **"Speech-to-Text with End-to-End Deep Learning":** Build a model that takes short audio clips and outputs text transcripts directly (use a dataset like LibriSpeech or Common Voice). Experiment by training both a traditional pipeline (e.g., MFCC features + classifier) and an end-to-end neural model, and compare their results to see where end-to-end methods shine—or struggle!
- 

## Whether to use End-to-end Deep Learning

### Pros and cons of end-to-end deep learning

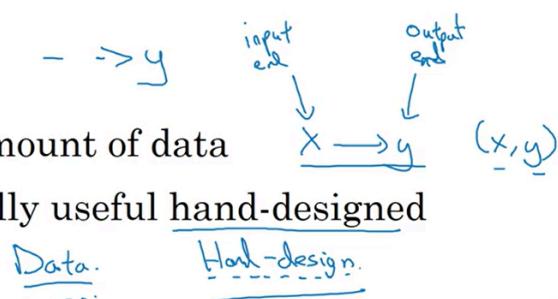
#### Pros:

- Let the data speak  $x \rightarrow y$
- Less hand-designing of components needed



#### Cons:

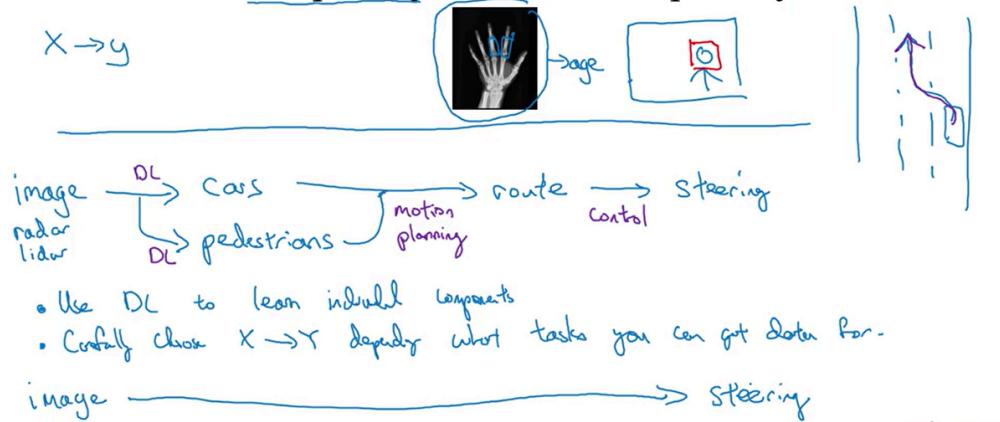
- May need large amount of data
- Excludes potentially useful hand-designed components



Andrew Ng

# Applying end-to-end deep learning

Key question: Do you have sufficient data to learn a function of the complexity needed to map  $x$  to  $y$ ?



## Core Concept:

End-to-end deep learning refers to training a neural network (or other machine learning model) to map inputs directly to outputs, learning the entire transformation from beginning to end, without splitting the problem into hand-designed intermediate steps. This approach lets data “speak for itself” by minimizing engineered features or manually crafted components.

## Breakdown:

### How It Works:

- The model receives raw input XXX and directly predicts output YYY.
- There are no intermediate representations or features designed by humans; the network is responsible for learning all relevant transformations and abstractions.
- The training process uses lots of labeled (X,Y)(X, Y)(X,Y) data, optimizing model parameters to minimize prediction error between predicted and actual YYY.
- By removing manually engineered subcomponents (like feature extractors), the model is flexible and potentially powerful but relies heavily on abundant data.

## Key Formulas:

A typical end-to-end deep learning mapping can be represented as:

$$Y = f_w(X) \quad Y = f_w(X)$$

$$Y = f_w(X)$$

Where:

- XXX: Raw input (e.g., audio, image)
- YYY: Desired output (e.g., transcription, label, action)
- $f_w$ : Deep neural network parameterized by weights  $w$
- Training seeks to minimize:

$$\min_w \sum_{i=1}^N L(f_w(X_i), Y_i)$$

$$w_{\text{mini}} = \frac{1}{N} \sum_{i=1}^N L(f_w(X_i), Y_i)$$

Where:

- LLL: Loss function (like MSE, cross-entropy)
- NNN: Number of training samples

*Intuitive Explanation:* The network learns any representation it needs, with no constraints except the input/output task and loss function.

---

## Why It's Used:

- *Problem Solved:* Reduces the need for manual design of features/components, allowing models to solve complex mapping tasks that may be challenging to decompose.
  - *Advantages:* Lets data determine the solution, which can uncover better or novel representations compared to hand-tuned pipelines. Especially powerful when vast labeled data is available.
- 

## Practical Insights:

### Real-World Applications:

- **Speech Recognition:** Modern voice-to-text systems (e.g., Google Voice, Siri) use deep networks trained end-to-end from audio waveform to transcribed text.
- **Image Classification:** Systems like Google Photos organize and classify images directly from pixels.

- **Autonomous Vehicles:** Some experimental approaches map raw sensor input (image, radar) directly to steering commands, although, in practice, hybrid systems with traditional pipelines are typically used.

### **Common Pitfalls:**

- **Requires Large Data:** End-to-end models usually need massive labeled datasets. Small datasets make it hard for the model to learn the necessary transforms.
- **Can Ignore Prior Knowledge:** By not using useful hand-designed components, models may miss out on proven, domain-specific techniques.
- **Overfitting:** Without enough data, the network may learn spurious correlations that don't generalize.
- **Debugging:** End-to-end systems, lacking interpretable steps, can be hard to troubleshoot.

### *How to Avoid Issues:*

- Use end-to-end learning only when ample, labeled data is available for the full task.
- For limited data, consider hybrid approaches: combine data-driven learning with human-engineered components or features.
- Regularization, data augmentation, and cross-validation can help with overfitting.

### **Creative Extension:**

#### **Project Idea:**

*Build an end-to-end handwritten digit recognition system.*

- **Task:** Train a deep neural network to map raw handwritten digit images directly to digit labels (0-9), using the MNIST dataset.
- **Goal:** Do not use any engineered preprocessing or feature extraction; rely solely on the raw image as input and optimize the network to predict the correct digit.
- **Extensions:** Try with both small and large subsets of the dataset to see the effect of data quantity on end-to-end performance. Explore adding/removing regularization or data augmentation and no

# Ruslan Salakhutdinov Interview

## 1. Personal Journey into Deep Learning

- Originally worked in the financial sector before considering a PhD.
- A chance encounter with Geoff Hinton led him into deep learning.
- Started a PhD with Hinton, coinciding with early work on Restricted Boltzmann Machines (RBMs), unsupervised pre-training, and deep learning (mid-2000s).

## 2. Early Research & Breakthroughs

- Co-authored a seminal paper on RBMs, enabling the resurgence of neural networks and deep learning.
- First experiments focused on autoencoders and using non-linear extensions of PCA.
- Early excitement working with datasets like MNIST digits and Olivetti faces for data compression and representation learning.
- Achieved strong results with deep autoencoders at a time when traditional methods couldn't train such deep models.

## 3. Evolution of Boltzmann Machines

- RBMs and stacked RBMs were crucial for early multilayer neural network training.
- Theoretical justification: Improved variational bounds as more layers are added.
- With the advent of GPUs (around 2009-2010), standard backpropagation (without pre-training) became possible, faster, and often more effective.
- Training Boltzmann Machines remains challenging: Markov Chain Monte Carlo and variational learning are less scalable than backpropagation.

## 4. Generative vs. Supervised Learning: Evolution of Thought

- Early efforts in deep learning centered on unsupervised and generative approaches because they solved the problem of working with mostly unlabeled data.
- Today, generative modeling (GANs, VAEs, deep energy models) is still a frontier, but supervised learning has driven most recent practical progress.

- There is still untapped potential in leveraging vast unlabeled data.

## 5. Advice for Aspiring Researchers and Practitioners

- Don't be afraid to try new or hard things, even if theory lags behind practical progress.
- "Just jump in"—experiment, try different things, embrace challenging problems.
- Don't rely only on high-level frameworks; understand the fundamentals (e.g., code backpropagation manually for deeper insight).
- Both academic and industrial research paths are viable; academia allows more freedom, while industry offers resources and the chance to impact millions.

## 6. Current Exciting Frontiers in Deep Learning

- **Unsupervised & semi-supervised learning:** Ability to utilize large unlabeled datasets.
- **Deep reinforcement learning:** Training agents in virtual environments and scaling algorithms for agent communication and interaction.
- **Reasoning and natural language understanding:** Building systems that can reason, have dialog, and answer questions intelligently.
- **Few-shot/one-shot/transfer learning:** Learning new tasks from minimal data, focusing on developing more human-like learning abilities.

## Key Takeaways to Remember

- Luck, openness, and mentorship can dramatically shape a research career.
- Compute advances have shifted the methods used in deep learning.
- Hands-on understanding of algorithms is invaluable.
- Both academia and industry play vital roles; choose based on personal preference for freedom vs. impact.
- The future of AI is likely in broader, more efficient use of unlabeled data, reinforcement learning, better reasoning, and more flexible learning approaches akin to human intelligence.[coursera](#)