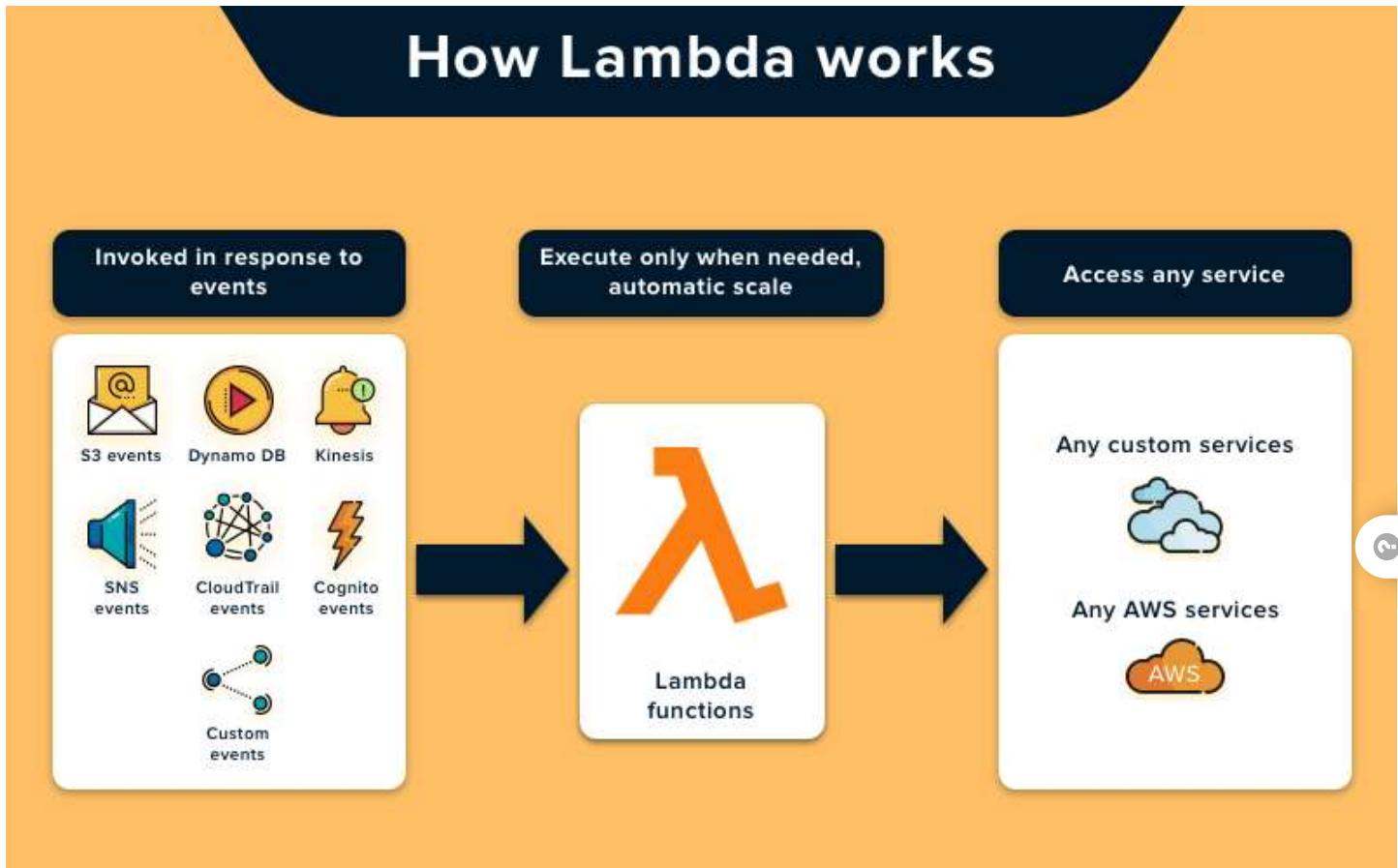


Session 4 - Basics of Neural Networks

Session 4: Basics of Neural Networks

- **AWS Lambda and Scheduler:** In-depth look into what AWS Lambda is, how to use it and Tie it up with Scheduler!
- **Understanding Data Types:** In-depth look at how images, text, and audio data are represented in computational systems.
- **Data Preprocessing Techniques:** Cleaning, normalizing, and preparing data for model training.
- **Data Augmentation:** Techniques to artificially expand visual, textual, or aural datasets.
- **Weights, Kernels, and Layers:** The building blocks of DNNs.
- **Types of Convolution Layers:** Convolution types, Strides, Multi-Channel Convolution
- **Embeddings:** Short introduction to Embeddings and Positional Embeddings
- **Activation Functions:** Different types and which ones to focus on.
- **CoreSets and Curriculum Learning:** What are they and how would they be useful for us (and NO we are not doing model distillation)



0) Prerequisites (once)

- **Google AI Studio API key:** create one and copy it somewhere safe.
- **AWS account:** if you don't have one, create at <https://aws.amazon.com/>
- [\(https://aws.amazon.com/\)](https://aws.amazon.com/) and sign in.

1) Log in to AWS & pick a region

1. Open <https://console.aws.amazon.com/>  (<https://console.aws.amazon.com/>) and Sign in.
2. Top-right, pick a region close to you (e.g., **Asia Pacific (Mumbai)**  **ap-south-1**). We'll use this same region everywhere below.

2) Store your Gemini key in Secrets Manager

1. In the AWS Console search bar, type **Secrets Manager** → open it.
2. Click **Store a new secret**.
3. **Secret type:** *Other type of secret*.
4. **Secret value:** choose **Plaintext** and paste your Gemini API key **as a single string** (no JSON).
 - Example (don't use this fake key):

AIzaSy...YourRealGeminiKey...123

5. Click **Next**.
6. **Secret name:**  **(exactly this)**.
7. Click **Next** → **Next** → **Store**.
8. Open the secret you just created and copy its **Secret ARN** (we'll need it when giving permissions).

3) Create the Lambda function (Python 3.12)

1. Go to **Lambda** → **Create function**.
2. **Author from scratch:**

- Name: **gemini-ping**
- Runtime: **Python 3.12**
- Architecture: **x86_64** (arm64 also fine)
- Permissions: **Create a new role with basic Lambda permissions**
- **Do not** attach a VPC (keep default “No VPC”) so it has internet access.

3. Click **Create function**.

4. In the function’s **Configuration** → **General configuration**, click **Edit**:

- **Memory**: 256 MB
- **Timeout**: 20 seconds Save.

4) Give Lambda permission to read the secret

1. In the Lambda function page, **Configuration** → **Permissions**.
2. Click the **Role name** link (opens the IAM role).
3. **Add permissions** → **Create inline policy**.
4. Choose the **JSON** tab and paste this (replace the **Resource** ARN with your secret’s ARN):

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "ReadGeminiSecret",
      "Effect": "Allow",
      "Action": ["secretsmanager:GetSecretValue"],
      "Resource": "arn:aws:secretsmanager:ap-south-1:YOUR_ACCOUNT_ID:secret:prod/gemini/api_key-*"
    }
  ]
}
```

1. **Review policy**, name it **AllowReadGeminiSecret**, **Create policy**.

5) Install the Gemini SDK as a Layer (one-time)

We'll build a Lambda Layer in **CloudShell** so dependencies match Lambda's Linux environment.

Click the **CloudShell** icon (top bar in the console). Wait for the terminal.

CloudShell's default is Python 3.9. We can still build a **Python 3.12–compatible** Lambda layer from CloudShell by **downloading cp312 wheels** and unpacking them. Do this exactly:

Build a Py3.12 x86_64 layer from CloudShell (works even if CloudShell is 3.9)

1. Prep folders and upgrade pip

```
python3 -m pip install --upgrade pip  
mkdir -p genai-layer/python/wheels
```

1. Download the correct wheels for CPython 3.12 / manylinux2014_x86_64

```
python3 -m pip download \  
--only-binary=:all: \  
--platform manylinux2014_x86_64 \  
--implementation cp \  
--python-version 312 \  
-d wheels \  
"google-genai==0.6.0" "pydantic==2.8.2" "pydantic-core==2.20.1"
```

1. Verify wheels exist

```
ls wheels | sed -n '1,999p'
```

You should see files like:

```
pydantic_core-2.20.1-cp312-...manylinux2014_x86_64.whl  
pydantic-2.8.2-py3-none-any.whl  
google_genai-0.6.0-py3-none-any.whl  
... (requests, google-auth, etc.)
```

1. Unzip all wheels into the layer folder

```
mkdir -p genai-layer/python
for f in wheels/*.whl; do
    unzip -o "$f" -d genai-layer/python > /dev/null
done
```

1. Zip the layer and publish

```
cd genai-layer
zip -r genai-layer.zip python
aws lambda publish-layer-version \
--layer-name google-genai-py312 \
--region <YOUR-REGION> \
--description "Google GenAI SDK layer (cp312, manylinux2014_x86_64)" \
--zip-file fileb://genai-layer.zip \
--compatible-runtimes python3.12
```

Copy the **LayerVersionArn** from the output.

1. Attach the layer and set runtime

2. Lambda → your function → Configuration → Runtime settings:

- **Runtime:** Python 3.12
- **Architecture:** x86_64

3. Lambda → Layers → Add a layer → Specify an ARN → paste the **LayerVersionArn** → Add.

4. Ensure only this py312 layer is attached.

5. Click **Deploy** in the Code tab (forces a fresh load).

6. Back to the AWS Console → Lambda → Layers → Create layer:

- Name: **google-genai-py312**
- Upload: the **genai-layer.zip** you just created (download from CloudShell first if the UI asks, or use the “Upload” from your local machine after downloading).
- Compatible runtimes: **Python 3.12**
- **Create** the layer.

7. Open your **gemini-ping** function → Layers → Add a layer → Custom layers → select **google-genai-py312** → Add.

6) Add environment variables

In your Lambda function: **Configuration** → **Environment variables** → **Edit** → **Add**:

- `GEMINI_SECRET_NAME` = `prod/gemini/api_key`
- `MODEL_NAME` = `gemini-2.5-flash`

Save.

7) Paste the Lambda code

In the function **Code** tab, ensure the **Handler** is `lambda_function.handler` (top-right “Runtime settings”). Then replace the default code with the following as `lambda_function.py` and **Deploy**:

```
import json
import os
import base64
import boto3
import logging

logger = logging.getLogger()
logger.setLevel(logging.INFO)

from google import genai # provided by the layer

SECRET_NAME = os.environ.get("GEMINI_SECRET_NAME", "prod/gemini/api_key")
MODEL_NAME = os.environ.get("MODEL_NAME", "gemini-2.5-flash")
REGION = os.environ.get("AWS_REGION") or os.environ.get("AWS_DEFAULT_REGION")

_API_KEY_CACHE = None
_CLIENT_CACHE = None
```

```

def _get_api_key():
    """Fetch and cache the API key from Secrets Manager."""
    global _API_KEY_CACHE
    if _API_KEY_CACHE:
        return _API_KEY_CACHE

    sm = boto3.client("secretsmanager", region_name=REGION)
    sec = sm.get_secret_value(SecretId=SECRET_NAME)
    s = sec["SecretString"]

    # Support either a bare string or a simple JSON structure
    try:
        data = json.loads(s)
        api_key = data.get("GEMINI_API_KEY") or next(iter(data.values()))
    except Exception:
        api_key = s

    _API_KEY_CACHE = api_key.strip()
    return _API_KEY_CACHE

```

```

def _get_client():
    global _CLIENT_CACHE
    if _CLIENT_CACHE:
        return _CLIENT_CACHE
    _CLIENT_CACHE = genai.Client(api_key=_get_api_key())
    return _CLIENT_CACHE

```

```

def _extract_prompt_from_event(event):
    """
    Supports:
    - Direct invoke: {"prompt": "..."}
    - Function URL GET: ?prompt=...
    - Function URL POST (JSON): {"prompt": ...}
    """

    if isinstance(event, dict):
        # direct JSON invoke
        p = event.get("prompt")
        if isinstance(p, str) and p.strip():
            return p.strip()

        # GET query
        q = event.get("queryStringParameters")

```

```

if isinstance(q, dict):
    p = q.get("prompt")
    if isinstance(p, str) and p.strip():
        return p.strip()

# POST body
body = event.get("body")
if body:
    if event.get("isBase64Encoded"):
        body = base64.b64decode(body).decode("utf-8", errors="ignore")
    try:
        data = json.loads(body)
        p = data.get("prompt")
        if isinstance(p, str) and p.strip():
            return p.strip()
    except Exception:
        pass

return "Explain why many teams moved from TensorFlow to PyTorch."

```

```

def handler(event, context):
    try:
        prompt = _extract_prompt_from_event(event)
        client = _get_client()

        resp = client.models.generate_content(model=MODEL_NAME, contents=prompt)
        text = getattr(resp, "text", None) or str(resp)

        # ---- NEW: log a compact JSON record (truncate for safety) ----
        out = {
            "ts": getattr(context, "aws_request_id", ""),
            "model": MODEL_NAME,
            "prompt": prompt,
            "response_preview": text[:2000],    # avoid 256 KB log limit
            "response_len": len(text),
        }
        logger.info(json.dumps(out, ensure_ascii=False))

        return {
            "statusCode": 200,
            "headers": {
                "content-type": "application/json",
                "Access-Control-Allow-Origin": "*",
                "Access-Control-Allow-Headers": "*",
            }
        }

    except Exception as e:
        logger.error(f"Error processing event: {e}")
        return {
            "statusCode": 500,
            "body": "Internal Server Error"
        }

```

```
        "Access-Control-Allow-Methods": "GET,POST,OPTIONS",
    },
    "body": json.dumps({"model": MODEL_NAME, "prompt": prompt, "response": text}, ensure_ascii=False),
}

except Exception as e:
    logger.exception("Lambda error")
    return {"statusCode": 500, "headers": {"content-type": "application/json"}, "body": json.dumps({"error": str(e)})}
```

8) Create a Function URL (so you can “see it”)

1. In your Lambda, go to **Configuration** → **Function URL** → **Create function URL**.
2. **Auth type: NONE** (public; fine for demo).
3. **Save**. Copy the **Function URL** shown at the top.

Test in a browser (GET):

```
https://<your-id>.lambda-url.<region>.on.aws/?prompt=One%20line%20on%20transformers
```

Test with cURL (POST JSON):

```
curl -X POST "<YOUR_FUNCTION_URL>" \
-H "content-type: application/json" \
-d '{"prompt":"Explain the attention mechanism in one paragraph."}'
```

You'll get a JSON response containing the Gemini answer.

9) (Optional) Console “Test” button

You can also test inside the Lambda console:

- Click **Test** → **Configure test event** → **Create new test event with:**

```
{ "prompt": "Explain why nearly everyone moved from TensorFlow to PyTorch?" }
```

- Click **Test**. You should see a **200** with the model's text.

10) Done — how “open/shutdown” works

Lambda automatically spins up, runs, and shuts down your function container. You don't need to close anything manually. Logs are in **CloudWatch Logs** (Lambda → Monitor → Logs).

Troubleshooting (quick)

- **No module named 'google'** → You forgot to attach the **Layer** to the function, or the layer zip didn't contain a top-level `python/` folder. Recreate the layer exactly as in Step 5.
- **AccessDeniedException for Secrets Manager** → The IAM inline policy's **Resource ARN** doesn't match your secret. Use the exact ARN from the secret page.
- **Timeouts** → Increase Lambda timeout to 20–30s (Step 3).
- **403 on Function URL** → Ensure Function URL **Auth type = NONE**.
- **No internet** → Don't place Lambda in a **VPC** unless you know how to give it a NAT gateway.

That's it? No

Next we'll use the **new EventBridge Scheduler** (recommended) to call your Lambda every 3 minutes. This gives you a super clear **Start (Enable) / Stop (Disable)** switch.

Follow these steps exactly in the same region as your Lambda.

Make a 3-minute scheduler (one time)

1. In the AWS console search bar, open **EventBridge**.
2. Left sidebar → **Schedules** → **Create schedule**.
3. **Name:** `gemini-ping-3m` (any name is fine). Optionally add a description.

4. Schedule pattern:

- Choose a **Recurring schedule**
- **Rate-based schedule** → Every `3` minutes
- **Flexible time window:** `Off` (keeps timing tight)
- **Start time:** *Run upon schedule creation* (default)
- **End time:** `None` (leave blank)
- Time zone: leave default (UTC) unless you prefer local.

5. Target details:

- **Target type:** **AWS Lambda**
- **Lambda function:** select your `gemini-ping` function
- **Execution role for your schedule:** choose **Create a new role for this schedule** (simplest & correct).
- **Invoke settings** → **Payload:** select **Constant (JSON)** and paste, for example:

```
{ "prompt": "Quick 1-line status: attention vs RNNs." }
```

(You can also leave `{}` because your handler has a default prompt.)

- Leave retries/dead-letter queue at defaults for now.

6. Next → **Create schedule**.

That's it. The schedule is **Enabled** by default, which means your “service” is **started** and will invoke the Lambda every ~3 minutes.

How to start and stop it later

- **Stop (Disable):** EventBridge → **Schedules** → click **gemini-ping-3m** → **Disable**.
- **Start (Enable):** same page → **Enable**.
- **Delete** removes it permanently.

Think of **Enable** = **start**, **Disable** = **stop**.

Verify it's running

- Lambda → **Monitor** → **Logs** → open the newest log stream; you should see one invocation about every 3 minutes.
- Or EventBridge → Schedules → **gemini-ping-3m** → **Monitoring** to see recent runs.

Change the prompt later (optional)

EventBridge → **Schedules** → **gemini-ping-3m** → **Edit** → under **Target details** change the **Payload JSON** → **Save**. Your Lambda will start receiving the new prompt next run.

Notes

- This path **invokes the Lambda directly** (not the Function URL), which your handler already supports.
- Costs are minimal: one invoke every 3 minutes is ~14,400/month; you remain within the free request tier in most accounts—compute time is what matters most.
- If you ever want to hit the **Function URL** instead (HTTP POST/GET), EventBridge Scheduler can also target an **HTTP endpoint**—happy to show those steps when you want.

Understanding Data Types

Images

- Representation:
 - Images are represented as multi-dimensional arrays (tensors) (can be JPG, PNG, BMP, etc)
 - But there are many other data formats, like DICOM Format (Medical):

Modality	Typical Channels	Description
X-ray	1 (Grayscale)	Single-channel grayscale images representing tissue density.
CT (Computed Tomography)	1 (Grayscale)	Multiple grayscale slices that form a 3D volume.
MRI (Magnetic Resonance Imaging)	1 (Grayscale)	Grayscale slices, representing tissue contrasts like T1, T2, etc.
PET (Positron Emission Tomography)	1 (Grayscale)	Grayscale images showing the distribution of a radioactive tracer.
Ultrasound	1 (Grayscale)	Grayscale imaging for soft tissue visualization.
Doppler Ultrasound	3 (RGB)	Color images representing blood flow velocity (Doppler).
Pathology/Microscopy	3 (RGB)	RGB images for tissue samples, often used in histology.
Multispectral Imaging	4-10+ (Multispectral)	Captures specific wavelengths across the spectrum, used for advanced diagnostics.

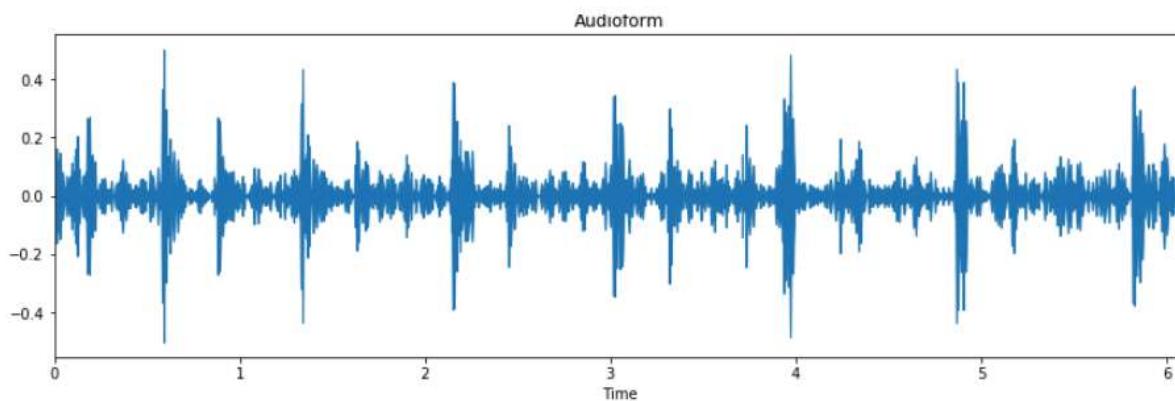
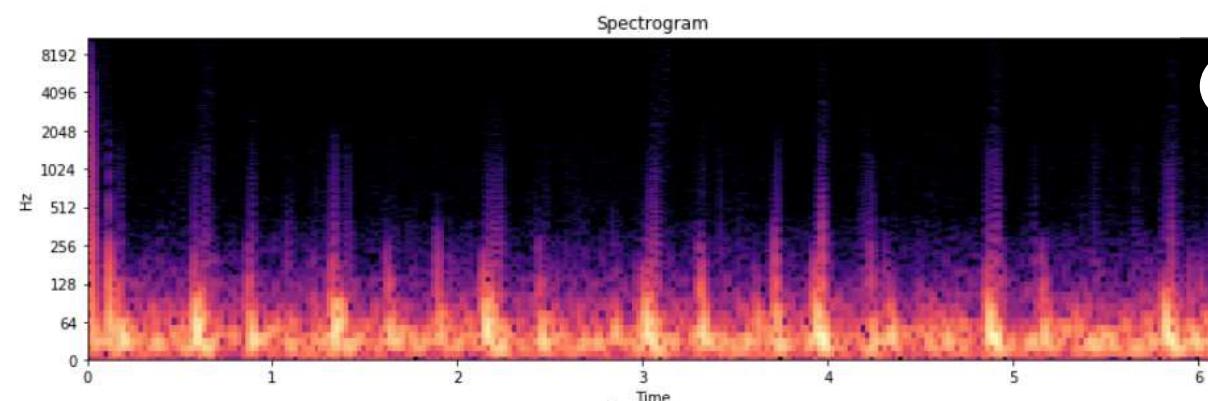
Modality	Typical Channels	Description
Hyperspectral Imaging	100-400+ (Hyperspectral)	Captures hundreds of narrow spectral bands, useful for detailed tissue analysis.
fMRI (Functional MRI)	4+ (Multichannel)	Multiple channels capturing different functional aspects like blood oxygenation.
PET-MRI or PET-CT	2 (Multimodal)	Combines grayscale data from PET and MRI/CT for a detailed diagnosis.
3D Ultrasound	Multiple slices (1 channel each)	Multiple grayscale images to form a 3D model.

Text

- Representation:
 - Text is represented as sequences of characters or [tokens](https://platform.openai.com/tokenizer) ↗ (<https://platform.openai.com/tokenizer>)
 - Tokenization: splitting text into meaningful units (words, subwords).
- Common Format: Plain text files, CSV, JSON, etc.
- Encoding: ASCII, UTF-8 (internal representation)
 - ASCII uses 7-bit encoding and supports 128 characters only (English letters, digits, punctuation)
 - UTF supports over 1.1 million characters, covering all languages and symbols
 - ASCII uses 1 byte per character, while UTF-8 supports 1 to 4 bytes depending on the character
 - UTF-8 is backward compatible with ASCII, making it efficient and widely used on the web
 - ASCII is limited to basic characters, whereas UTF can represent virtually any character from any language

Audio

- Representation:

- Audio can be shared as either 1D or 2D data, depending on how it is represented:
 - 1D Data (Time-domain Representation):
 - Raw audio signals (like a waveform) are stored as 1D data, where the amplitude of the sound is recorded at regular intervals. This is what you typically get in a sound file like .wav or .mp3. The data consists of a sequence of amplitude values over time. Later we would use FTs to extract individual frequencies:
- 
- 
- 2D Data (Frequency-Domain Representation):
 - When audio is transformed using techniques like STFT it becomes 2D data, which is often visualized as a spectrogram. The two dimensions are Time and Frequency. The original raw audio can be converted into 2D, and how we can see each frequency change over time:
- Sampling rate: The number of samples per second (e.g. 44.1kHz)
 - Common Formats: WAV, MP3, FLAC, etc

3D Data:

- Representation:

- 3D data is represented using geometric constructs
- Points Clouds: Sets of points in 3D space
- Meshes: Vertices, edges, and faces defining surfaces
- Common Formats: OBJ, STL, PLY

Creating Datasets in PyTorch

```
from torch.utils.data import Dataset

class CustomDataset(Dataset):
    def __init__(self):
        # Initialization code
        pass

    def __len__(self):
        # Returns the total number of samples
        return 0

    def __getitem__(self, idx):
        # Generates one sample of data
        return None
```

Data Preprocessing Techniques

Images:

- Resizing: to ensure all images are the same size
- Normalization: Scaling pixel values to a standard range (normally between +-1)

```
transform = transforms.Compose([
    transforms.Resize((128, 128)),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.5], std=[0.5])
])
```

1. Resize: from anything to 128x128 (assume its image and data is saved in int8 (0-255))

2. **ToTensor**: a lot is happening here:

1. converts into a **torch.floatTensor**

2. **min-max normalization**:

1. if data in int8:

1. scales by 1/255 (min-max normalization), so data after this step is [0, 1]

2. if data in any other format:

1. doesn't scale

3. changes HxWxC to CxHxW (H: height, W: width, C: channel) as Pytorch expects "channels" first

3. **normalize**: (0.5 values mentioned here are [PIDOOMA](#) ↗

(<https://en.wiktionary.org/wiki/PIDOOMA>), and must actually be calculated (S5))

1. $x_{norm} = \frac{(x-\mu)}{\sigma}$ for each pixel in each channel

Text:

- Tokenization: Splitting text into tokens
- Padding/Truncating: To ensure sequences are of equal length
- Embedding: Mapping tokens to vectors

```
from torch.nn.utils.rnn import pad_sequence

def collate_fn(batch):
    texts, labels = zip(*batch)
    texts = [torch.tensor(text) for text in texts]
    texts_padded = pad_sequence(texts, batch_first=True, padding_value=vocab['<pad>'])
    labels = torch.tensor(labels)
    return texts_padded, labels
```

Audio:

- Resampling: To standardize sampling rates
- Feature Extraction: Converting raw audio to features like MFCC

```
transformation = torchaudio.transforms.MFCC(sample_rate=16000, n_mfcc=40)
```

3D Data:

- Normalization: Scaling to a unit sphere
- Centering: Shifting to the origin

```
def normalize(vertices):
    centroid = vertices.mean(dim=0)
    vertices = vertices - centroid
    max_dist = torch.max(torch.norm(vertices, dim=1))
    vertices = vertices / max_dist
    return vertices
```

Data Augmentation

Images:

- Random Horizontal Flip
- Random Rotation
- Color Jitter
- [others ↗\(https://albumentations.ai/\)](https://albumentations.ai/)

```
augmentations = transforms.Compose([
    transforms.RandomHorizontalFlip(),
    transforms.RandomRotation(10),
    transforms.ColorJitter(brightness=0.2, contrast=0.2),
    transforms.ToTensor()
])
```

Why are augmentations applied before ToTensor and normalization?

Text:

- Synonym Replacement
- Random Insertion

```
import random
from nltk.corpus import wordnet

def synonym_replacement(sentence, n):
    words = sentence.split()
    for _ in range(n):
        word_to_replace = random.choice(words)
        synonyms = wordnet.synsets(word_to_replace)
        if synonyms:
            synonym = synonyms[0].lemmas()[0].name()
            words = [synonym if word == word_to_replace else word for
```

```
word in words]
    return ' '.join(words)
```

Audio:

- Time Stretching
- Pitch Shifting
- Adding Background Noise

```
def add_background_noise(signal, noise_factor=0.005):
    noise = torch.randn_like(signal)
    augmented_signal = signal + noise_factor * noise
    return augmented_signal
```

3D Data:

- Random Rotation
- Scaling
- Adding Noise

```
def random_rotation(vertices):
    theta = random.uniform(0, 2 * math.pi)
    rotation_matrix = torch.tensor([
        [math.cos(theta), -math.sin(theta), 0],
        [math.sin(theta), math.cos(theta), 0],
        [0, 0, 1]
    ])
    return vertices @ rotation_matrix
```



Model Architecture

```
import torch.nn as nn

class ImageClassifier(nn.Module):
    def __init__(self):
        super(ImageClassifier, self).__init__()
        self.conv_layers = nn.Sequential(
            nn.Conv2d(3, 32, kernel_size=3),
            nn.ReLU(),
            nn.MaxPool2d(2),
            # Add more layers as needed
        )
        self.fc_layers = nn.Sequential(
            nn.Linear(32 * 62 * 62, 128),
            nn.ReLU(),
            nn.Linear(128, 10) # For CIFAR-10
        )

    def forward(self, x):
        x = self.conv_layers(x)
        x = x.view(x.size(0), -1)
        x = self.fc_layers(x)
        return x
```

Training Loop:

```
model = ImageClassifier()
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)

for epoch in range(num_epochs):
    for images, labels in train_loader:
        outputs = model(images)
        loss = criterion(outputs, labels)
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```

Transformers for Text

Model Architecture

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class DecoderLayer(nn.Module):
    def __init__(self, d_model, num_heads, d_ff, dropout=0.1):
        super(DecoderLayer, self).__init__()
        self.self_attn = nn.MultiheadAttention(d_model, num_heads, dropout=dropout)
        self.norm1 = nn.LayerNorm(d_model)
        self.ffn = nn.Sequential(
            nn.Linear(d_model, d_ff),
            nn.ReLU(),
            nn.Linear(d_ff, d_model)
        )
        self.norm2 = nn.LayerNorm(d_model)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x):
        # Self-attention
        attn_output, _ = self.self_attn(x, x, x)
        x = self.norm1(x + self.dropout(attn_output))

        # Feed forward network
        ffn_output = self.ffn(x)
        x = self.norm2(x + self.dropout(ffn_output))

    return x

class TransformerDecoder(nn.Module):
    def __init__(self, vocab_size, d_model, num_heads, num_layers, d_ff, max_len, dropout=0.1):
        super(TransformerDecoder, self).__init__()
        self.embedding = nn.Embedding(vocab_size, d_model)
        self.pos_embedding = nn.Embedding(max_len, d_model)
        self.layers = nn.ModuleList([
            DecoderLayer(d_model, num_heads, d_ff, dropout) for _ in range(num_layers)
        ])
        self.fc_out = nn.Linear(d_model, vocab_size)
        self.dropout = nn.Dropout(dropout)

    def forward(self, x, pos):
        # Input embeddings + positional embeddings
        x = self.embedding(x) + self.pos_embedding(pos)
        x = self.dropout(x)

        # Pass through decoder layers
        for layer in self.layers:
            x = layer(x)

        # Output projection
        logits = self.fc_out(x)
        return logits

# Example usage
vocab_size = 10000
```

```
d_model = 512
num_heads = 8
num_layers = 6
d_ff = 2048
max_len = 100

model = TransformerDecoder(vocab_size, d_model, num_heads, num_layers, d_ff, max_len)
```

Training Loop:

```
from transformers import AdamW

optimizer = AdamW(model.parameters(), lr=2e-5)

for epoch in range(num_epochs):
    for texts, labels in train_loader:
        inputs = tokenizer(texts, return_tensors='pt', padding=True,
                           truncation=True)
        outputs = model(**inputs, labels=labels)
        loss = outputs.loss
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```

Building First Neural Networks

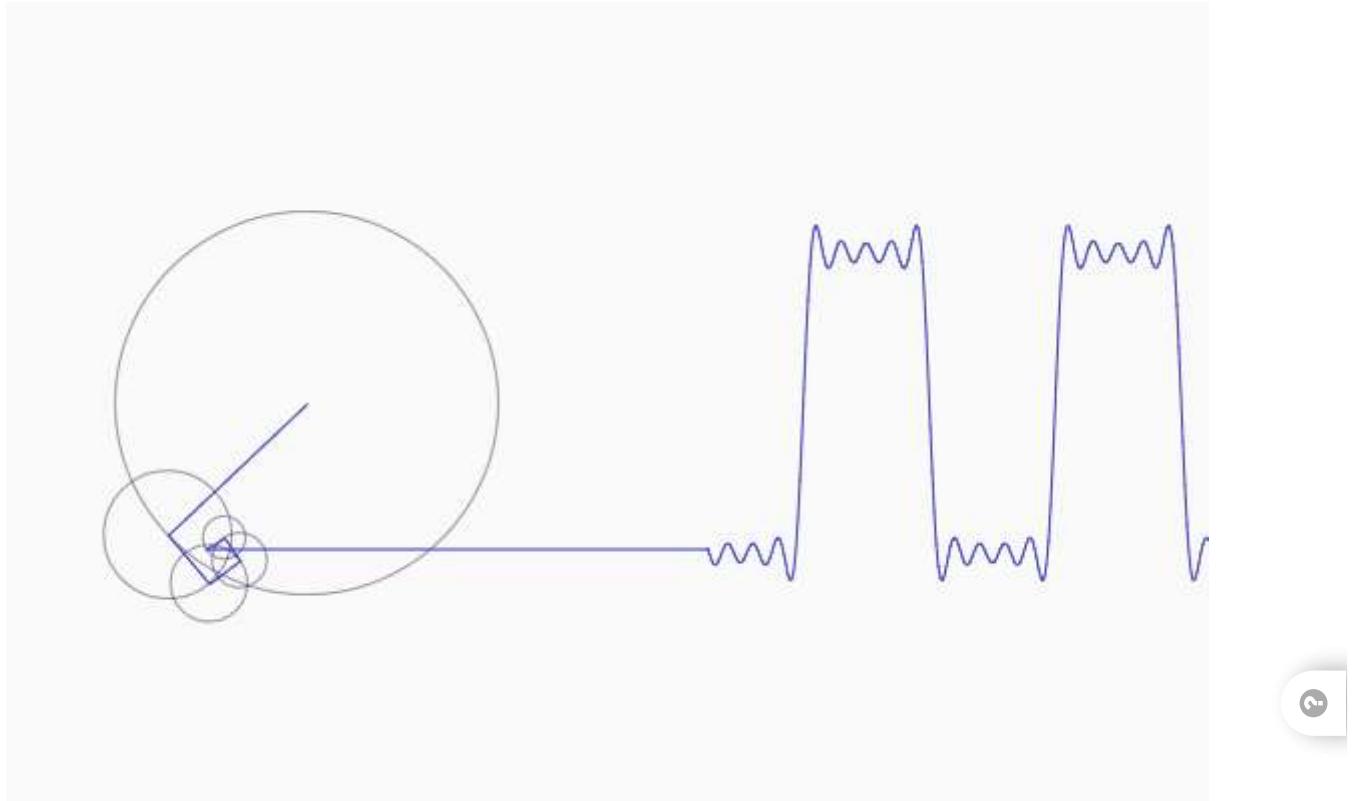
- **Weights, Kernels, Activation and Layers:** The building blocks and DNNs
- **Basics of Model Deployment:** Steps involved in deploying AI models to production.

- **Introduction to CI/CD Pipelines:** Automating testing and deployment using tools like GitHub Actions.

Let us first understand "why neural networks work so well?"

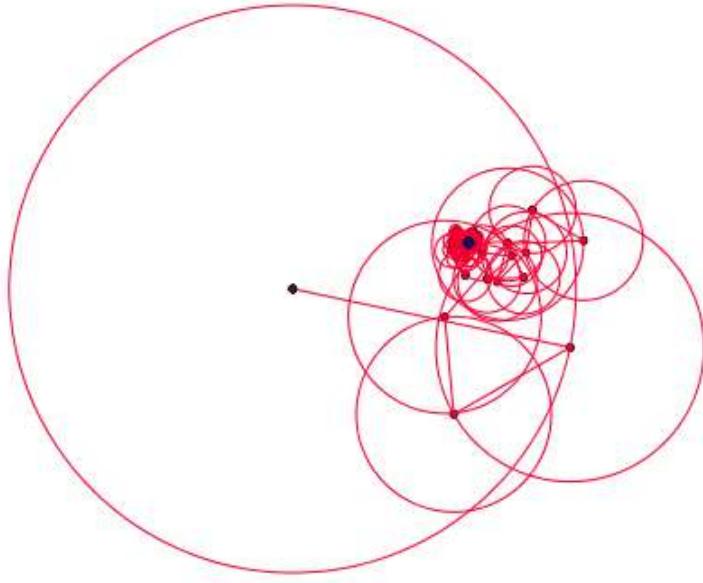
Most of the problems in the world are not simple, i.e. real-world "data" cannot be modeled with a single line. Play with this tool below:

Most of the problems are non-linear, and we need non-linear tools. Just observe how multiple Fourier transform can get together to create something (here a square wave):



Look at this amazing work by Jazzemon on [Fourier ↗ \(https://www.jezzamon.com/fourier/\)](https://www.jezzamon.com/fourier/)

And we can create complex things as well:



If you ever want to understand them deeply, just watch [this](#)  (<https://www.youtube.com/watch?v=spUNpyF58BY>) video (good for intuition)

Neural Networks can have a large number of free parameters (or weights between interconnected units), and this gives them the flexibility to fit highly complex data (when trained correctly) that other models are too simple to fit. **And Fourier Transformers can be considered to be a kind of neural network.** (If you want to understand this more, when you're free (like when taking ) ask ChatGPT about it).

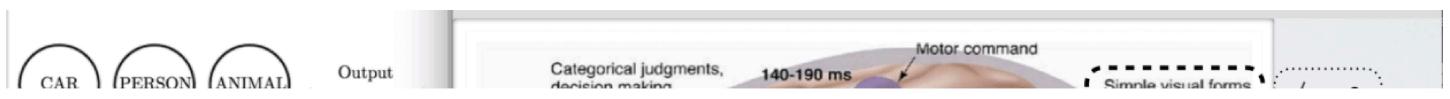
This model complexity brings with it the problems of training such a complex network and ensuring the resultant model generalizes to the examples it is trained on. Typically, neural networks [require](#)  (<https://stackoverflow.com/questions/38595451/why-do-neural-networks-work-so-well>) large volumes of training data, which other models don't.

What kind of data are we dealing with?

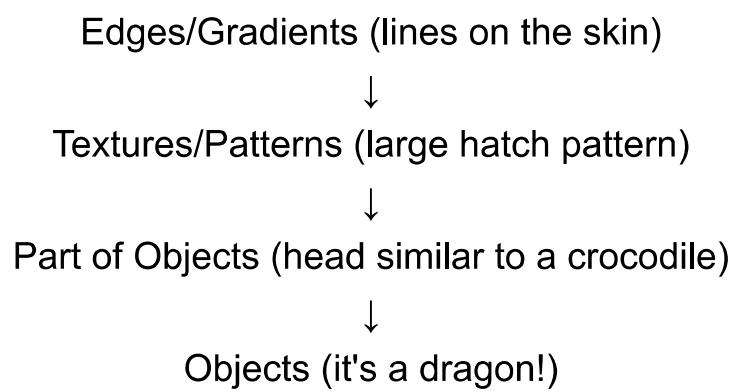
Below you see a car:



This image is broken down into smaller "features" by the 4 layers in our visual cortex:



Only after an image is broken down into its features, our brain combines it back into logical pieces:



Resources

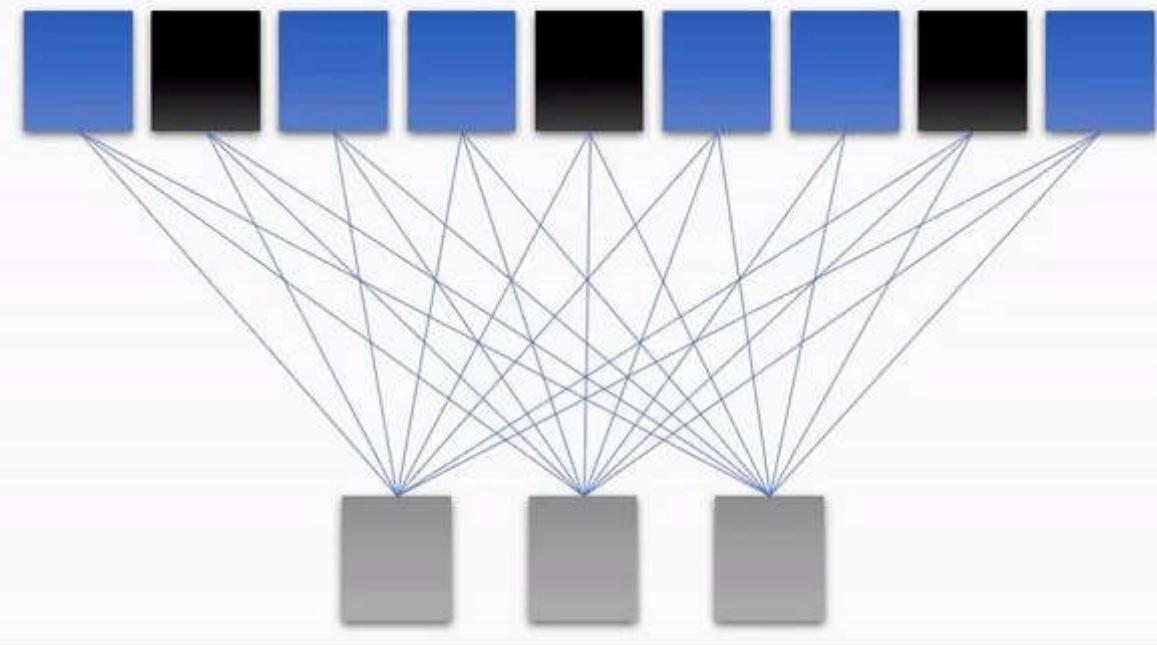
1× 0.5× 0.25×

Rerun



Feed-forward Networks

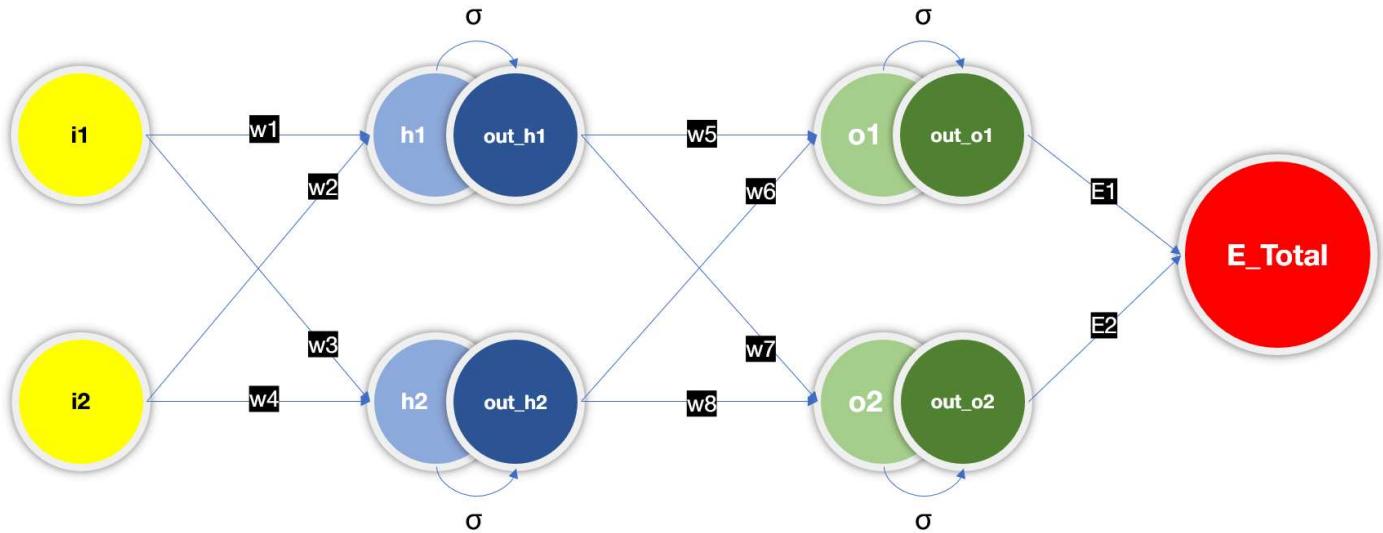
We discussed in the last session that when we use FC layers, we lose the temporal data. Look at that animation below:



In an FC layer, there is no concept of spatial or temporal information. Neurons are responsible for learning an abstract representation of the data. But what if??



Feedforward networks have the following characteristics



Neurons or perceptrons are arranged in a layer, with the first layer taking in inputs and the last layer producing the output. The middle layers have no connections with the external world and hence are called hidden layers. Each neuron or perceptron in one layer is connected to every perceptron on the next layer. Hence information is constantly "fed forward" from one layer to the next.

There are no connections among the perceptions in the same layer.

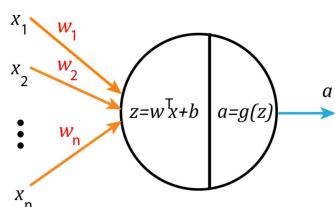
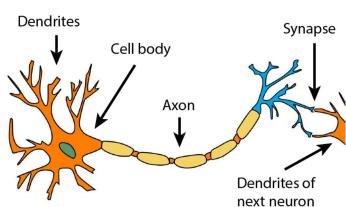
Each of the circles that you see is actually a neuron, but **each line** you see is the **weight** that we are training and is of importance to us. Those circles are "temporary" values that will be stored. Once you train the model, *lines are what all matter!*

These lines are denoted by **w**'s, where i and j are the neurons they connect together.



What are neurons?

A neuron is a fundamental unit of our brain. Neuron, w.r.t. our brain, consists of a small "memory storage" or a "signal", as well as a "small computational unit". When we refer to neurons in DNNs/NNs we only consider a small "memory storage" or a "signal" and keep the computation unit outside. This computation unit consists of two elements, a **weight**, and an **activation** function. Each neuron in both cases has input connections. Input connections to a brain neuron are called a dendrite and output connection is called an axon. Both are called just input and output weights in NNs.



Assuming these are **n** connections coming in, the output of the neuron (after using the weights and activation function) can be represented as:

$$a = \tanh(z) = \tanh(\sum_{i=1}^n (w_i + b_i))$$

where **b** is a bias with which we are stuck for historical reasons. **tanh** is playing the role of an activation function here.

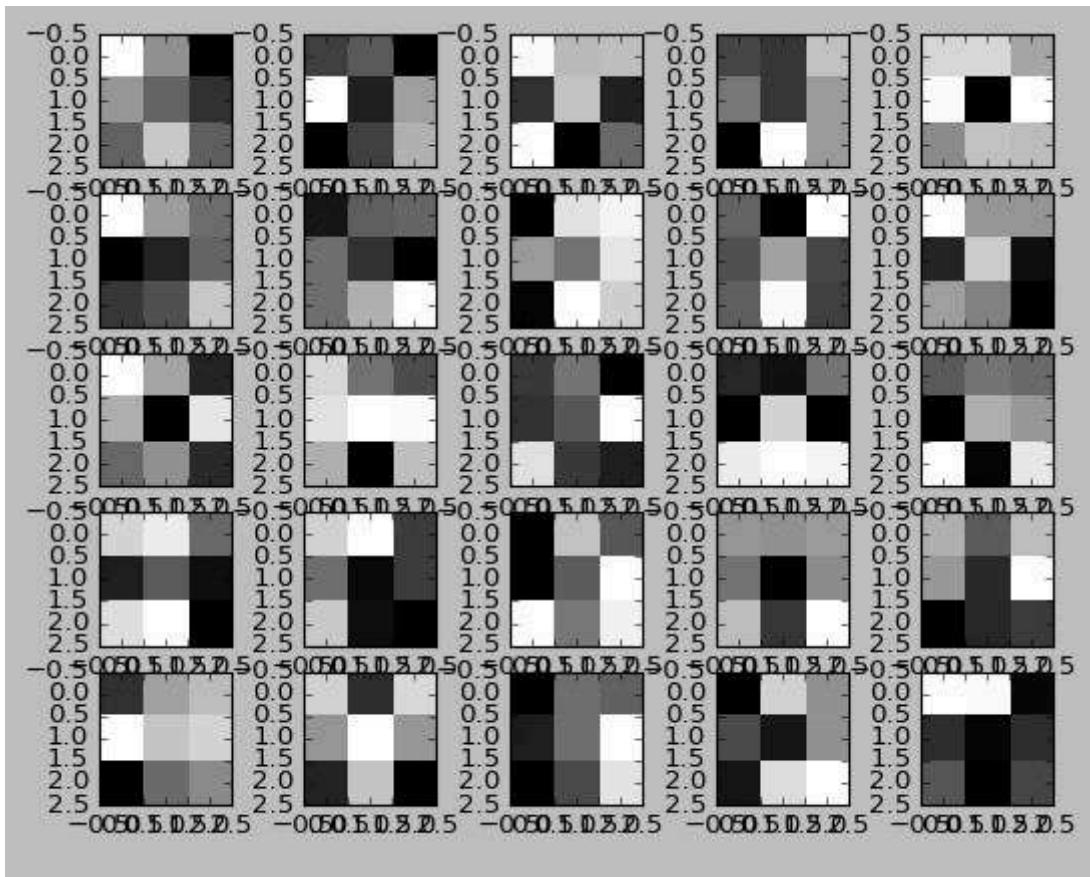
In our brain, we have different kinds of neurons doing different things (the brain's activation function) with the information coming in. In the case of NNs, we generally have a single activation like tanh/sigmoid/ReLU, etc.

Let's quickly build up an Excel sheet and make a Feedforward NN, that we will be using for performing backpropagation in the next session.

The neuron is a fundamental unit of our brain. Neuron, w.r.t. our brain consists of a small "memory storage" or a "signal", as well as a "small computational unit". When we refer to neurons in DNNs/NNs we only consider a small "memory storage" or a "signal" and keep the computation unit outside.

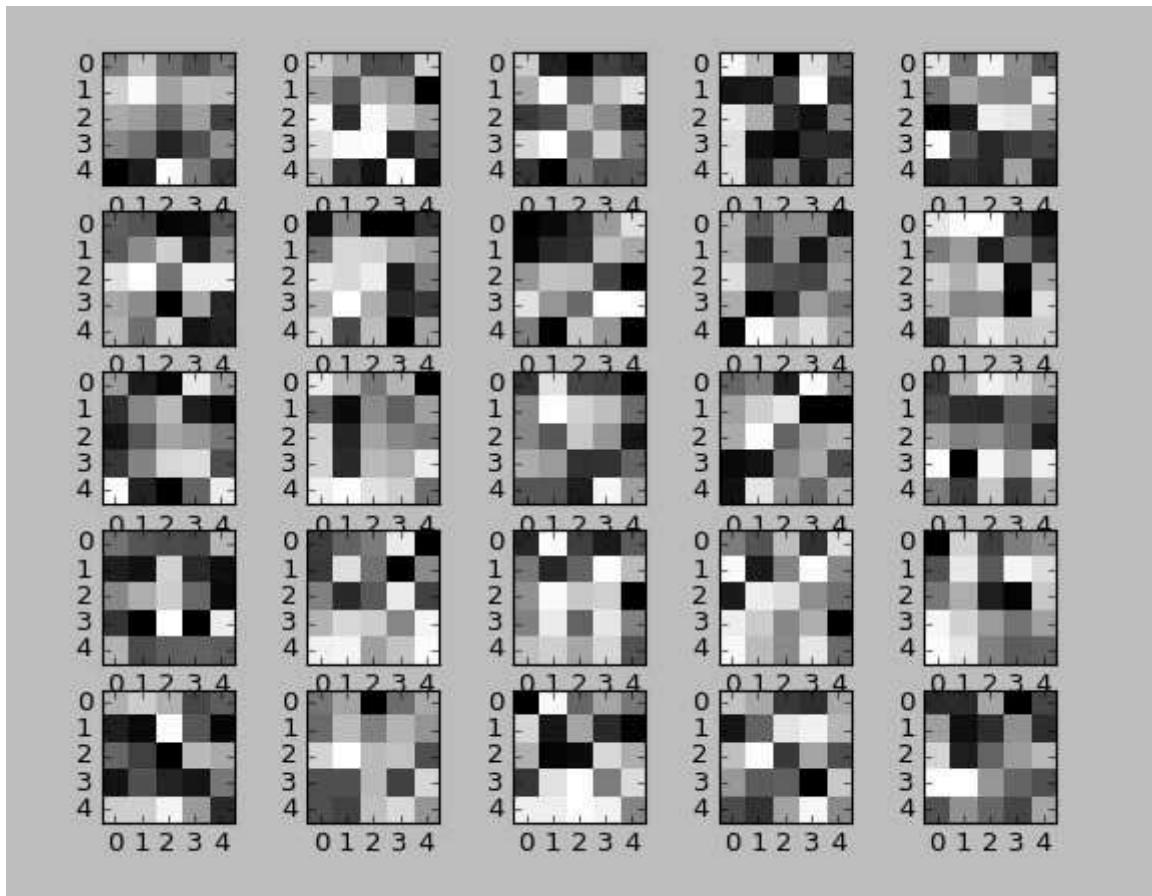
This computation unit consists of two elements: a **weight** and an **activation** function. Each neuron in both cases has input connections. Input connections to a brain neuron are called dendrites, and output connection is called axons. Both are called just input and output weights in NNs.

Let's look at what 3x3 kernels are **trying to extract**:

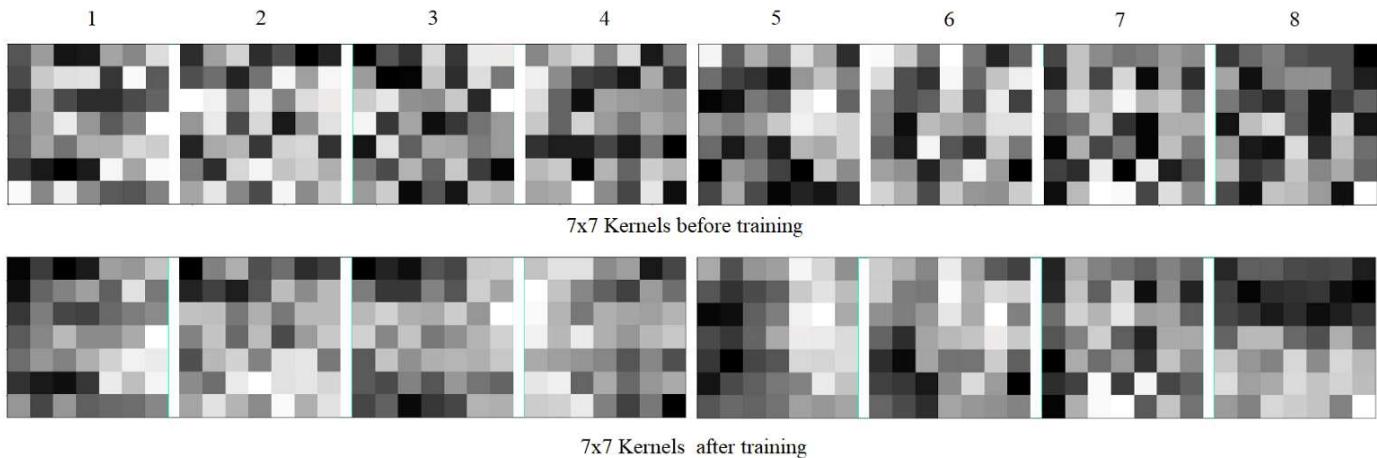


We really cannot find any pattern here. A 3x3 block is a very small area to gather any perceivable information (for human eyes), especially when we are referring to an image of size 400x400.

How about 5x5 kernels?

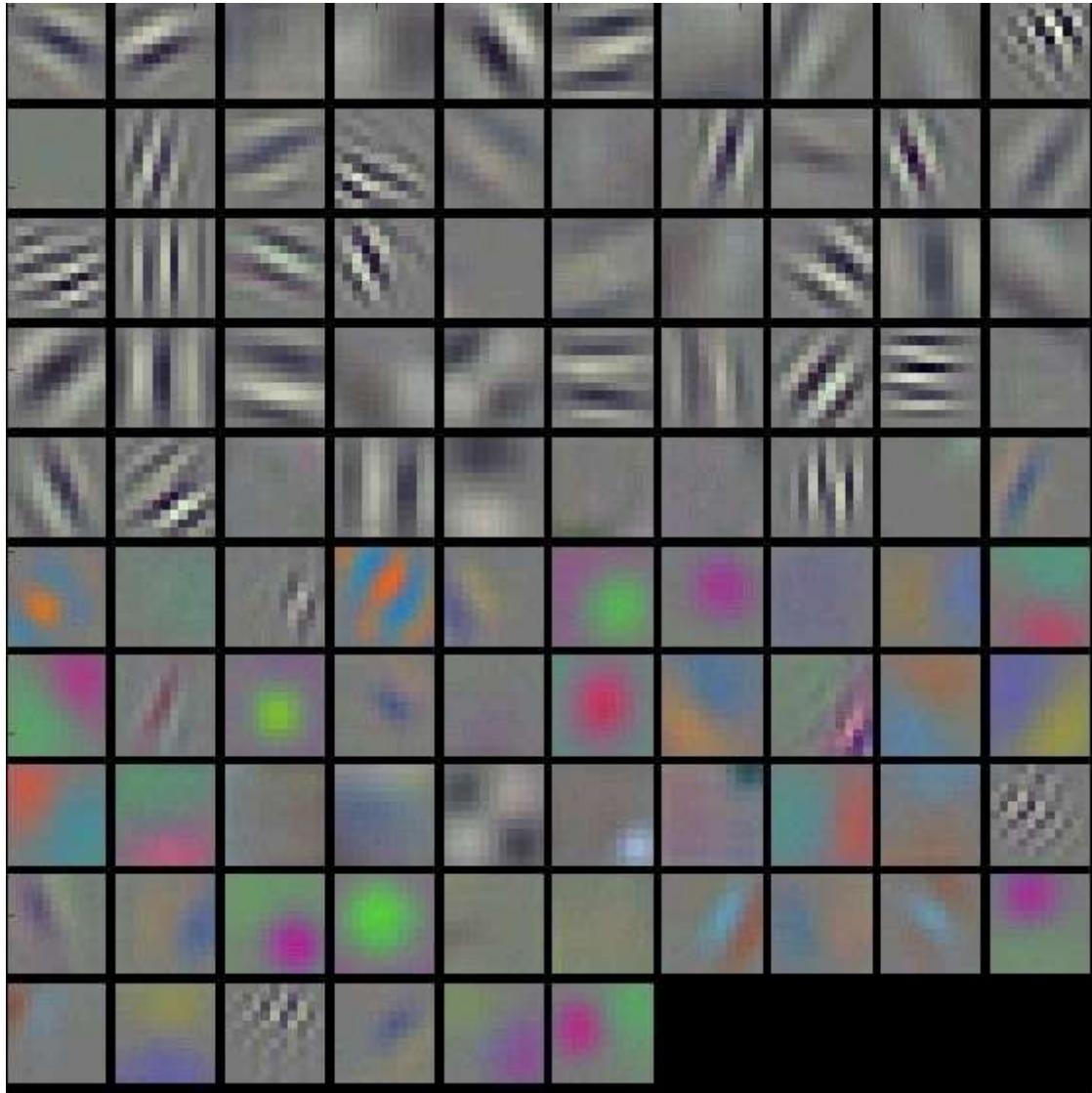


Even at 5x5 (from large image sizes) we cannot make out the patterns being extracted. In fact, at 7x7 we might also fail (look below):



11x11 - the turning points for +200 images

It is only at the receptive fields of around 11x11 when we would be able to make out patterns being extracted:

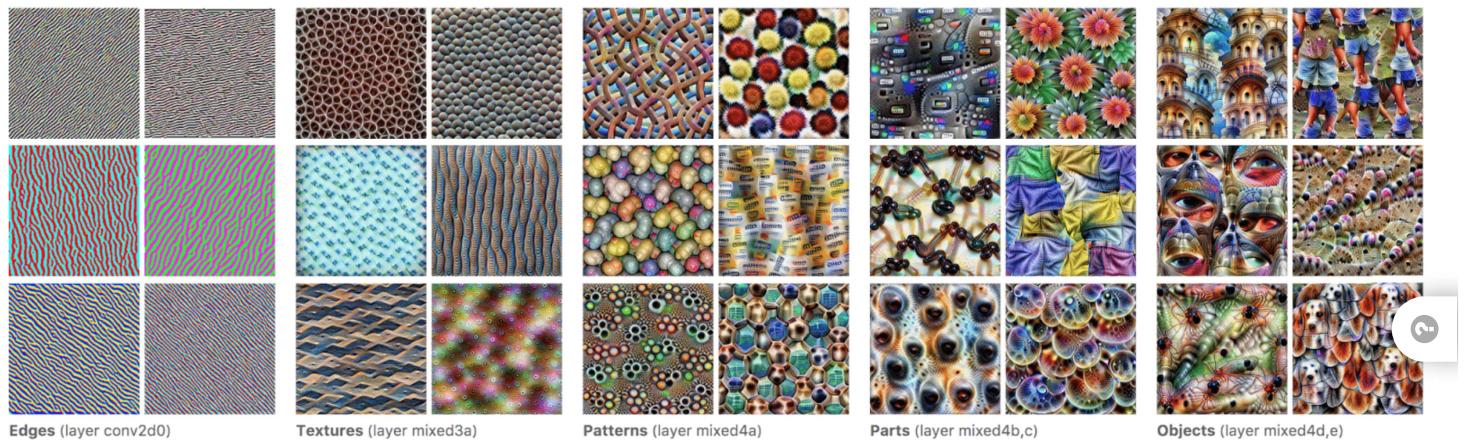


Let's visualize [\(https://ezyang.github.io/convolution-visualizer/\)](https://ezyang.github.io/convolution-visualizer/) how these things work [\(\(https://brilliantwavetech.com/cnn-visualization.html\)\)](https://brilliantwavetech.com/cnn-visualization.html): here [\(https://poloclub.github.io/cnn-explainer/\)](https://poloclub.github.io/cnn-explainer/), and here. [\(https://adamharley.com/nv_vis/\)](https://adamharley.com/nv_vis/)

4 Blocks - Edges/Gradients, Textures/Patterns, Part of Objects, and Objects.

All this while we have been discussing that we extract edges and gradients, textures, patterns, parts of objects, and then objects. But is it true?

It is indeed true! Look at the visualization below, where we can see what our kernels are extracting:



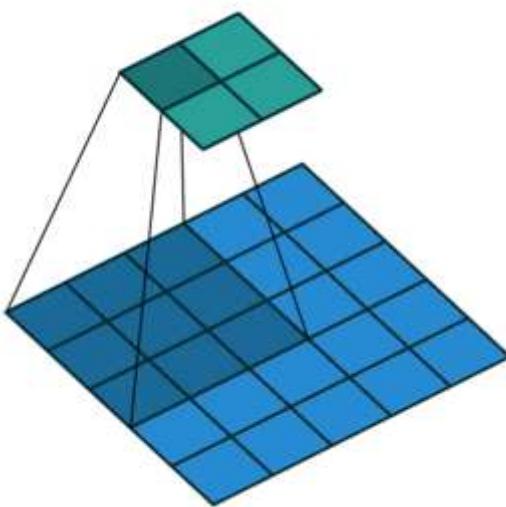
Some of the very important weblinks that we need to take a look at:

[CNN Explorer](https://poloclub.github.io/cnn-explainer/) ➔ (<https://poloclub.github.io/cnn-explainer/>)

[Feature Visualization](https://distill.pub/2017/feature-visualization/) ➔ (<https://distill.pub/2017/feature-visualization/>)

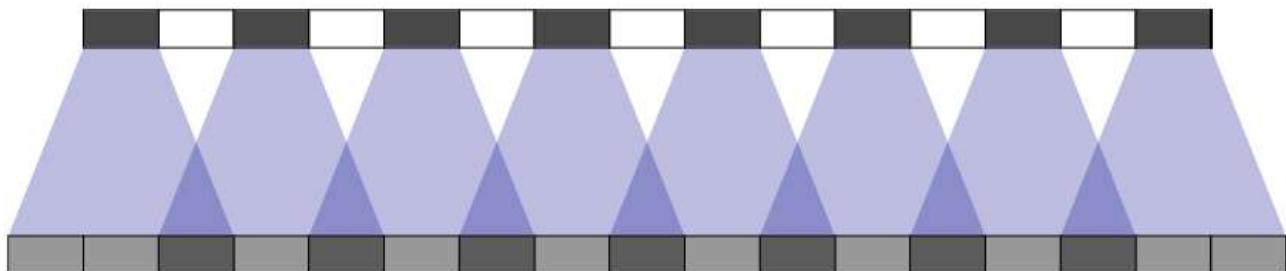
[Receptive Fields](https://distill.pub/2019/computing-receptive-fields/) ➔ (<https://distill.pub/2019/computing-receptive-fields/>)

Let us first discuss the consequences of using strides:

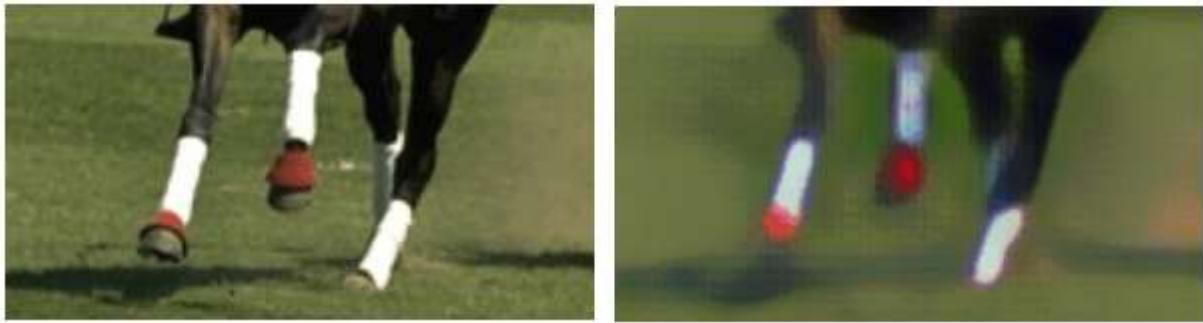


We learned in the last lecture that when we convolve with the standard way (stride of 1), we cover most of the pixels 9 times (we covered only 1 pixel 9 times on a channel of size 5x5, but as the channel size increases, we cover more and more pixels 9 times e.g. on a channel of size 400x400, 396x396 pixels would be covered 9 times).

But, when we convolve with a stride of more than 1, we would be covering some pixels more than once. And this is not good, as we are creating an island of extra information spread around in a repeating pattern. The image below would help:



As you can see, we are spreading the information unevenly. We are blurring the inputs as can be seen in the image below:



Let's check out [distill ↗ \(https://distill.pub/2016/deconv-checkerboard/\)](https://distill.pub/2016/deconv-checkerboard/) again.

We will come across this checkerboard issue again in super-resolution algorithms.

Observe the last image in the sequence of images below:

Perceptual Losses for Real-Time Style Transfer and Super-Resolution



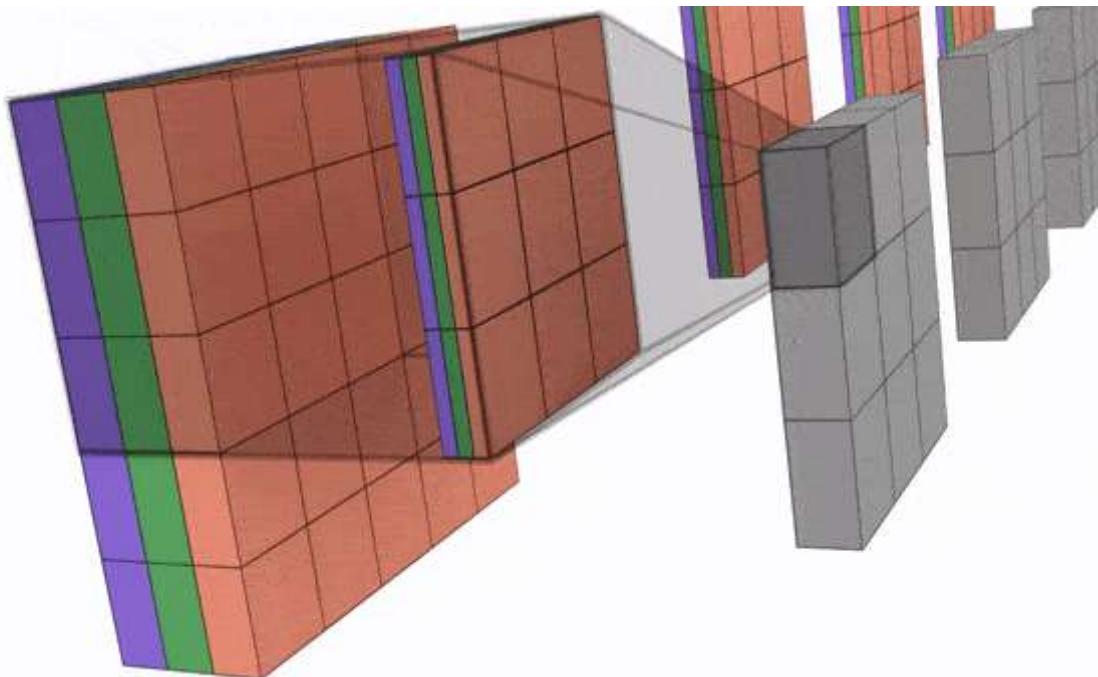
But then why do we use Convolutions with strides?

1. it reduces the amount of computation required by the network
2. it reduces the size of the output, hence reducing the overall number of layers to achieve the target RF
3. it increases the RF drastically
4. it adds additional invariant feature learning capabilities to the network
5. it improves the speed of the inferencing.

Listing two major disadvantages as well:

1. it causes checkerboard issues in the channels/images
2. it reduces the spatial resolution of the output, which means it makes it more difficult for further layers to learn detailed, fine-grained features.

Multi-Channel Convolution



It should be clear that we are using a $3 \times 3 \times 3 \times 4$ kernel above.

Another way of looking at them is this:

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+ 1 = -25

Bias = 1

-25					...
					...
					...
					...
...

Output

Let's build the first block and see why we are still missing a very strong Avenger from our team, Ant-Man!

The Rich-man's problem!

Let's build a DNN.

400x400x3 | (3x3x3)x32 | 398x398x32 RF of 3x3

398x398x32 | (3x3x32)x64 | 396x396x64 RF of 5X5

396x396x64 | (3x3x64)x128 | 394x394x128 RF of 7X7

394x394x128 | (3x3x128)x256 | 392x392x256 RF of 9X9

392x392x256 | (3x3x256)x512 | 390x390x512 RF of 11X11

MaxPooling

195x195x512 | (?x?x512)x32 | ?x?x32 RF of 22x22

.. 3x3x32x64 RF of 24x24

.. 3x3x64x128 RF of 26x26

.. 3x3x128x256 RF of 38x28

.. 3x3x256x512 RF of 30x30

Some points to consider before we proceed:

1. In the network above, the most important numbers for us are:
 1. **400x400**, as that defines where our edges and gradients would form
 2. **11x11**, as that's the receptive field that we are trying to achieve before we add transformations (like channel size reduction using MaxPooling)
 3. **512** kernels, as that is what we would need at a minimum to describe all the edges and gradients for the kind of images we are working with (ImageNet)
2. We have added 5 layers above, but that is **inconsequential**, as we aimed to reach the 11x11 receptive field. For some other datasets, we might have reached the required RF for edges&gradients, say after 4 or 3 layers
3. We are using 3x3 because of the benefits it provides; our ultimate aim is to reach the receptive field of 11x11
4. We are following 32, 64, 128, 256, and 512 kernels, but there are other possible patterns. We are choosing this specific one as this is expressive enough to build the 512 final kernels we need, and since this is an experiment, we could, later on, reduce the kernels depending on what hardware we pick for deployment.
5. The receptive field of 30x30 is again important for us because that is where we are "hoping" from textures.
6. The 5 Convolution Layers we see above form the "**Convolution Block**".

The question we left unanswered in the last session was, How should we reduce the number of kernels. We cannot just add 32, 3x3 kernels, as that would re-analyze all the 512 channels and give us 32 kernels. This is something we used to do before 2014, and it works, but intuitively, we need something better. We have 512 features now. Instead of evaluating these 512 kernels and coming out with 32 new ones, it makes sense to combine them to form 32 mixtures. That is where 1x1 convolution helps us.

Think about these few points:

1. Wouldn't it be better to merge our 512 kernels into 32 richer kernels that could extract multiple features that come together?
2. 3x3 is an expensive kernel, and we should be able to figure out something lighter, and less computationally expensive method.
3. Since we are merging 512 kernels into 32 complex ones, it would be great if we do not pick those features that are not required by the network to predict our images (like backgrounds).

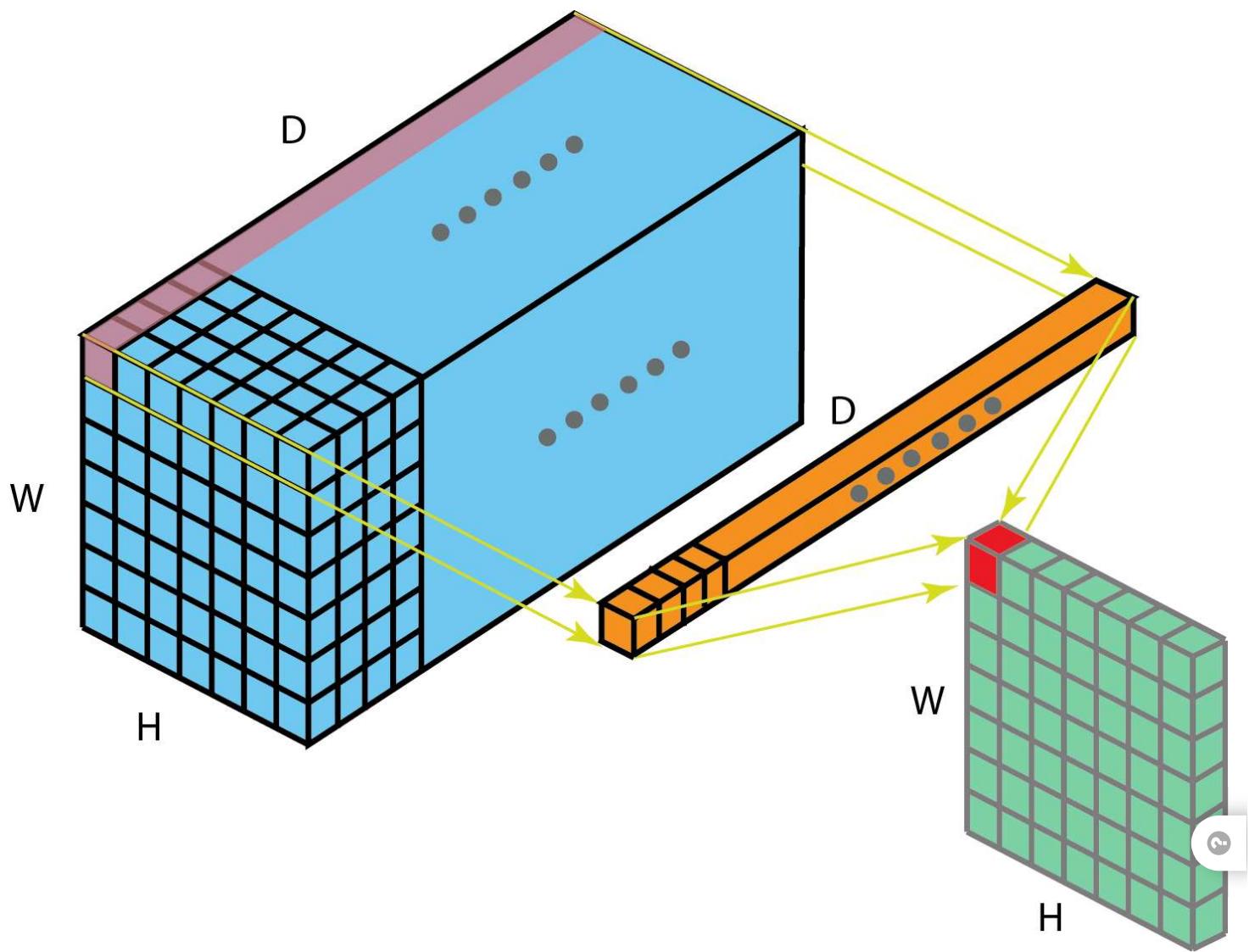
1x1 provides all these features.

1. 1x1 is computation less expensive.
2. 1x1 is not even a proper convolution, as we can, instead of convolving each pixel separately, multiply the whole channel with just 1 number
3. 1x1 is merging the pre-existing feature extractors, creating new ones, keeping in mind that those features are found together (like edges/gradients that make up an eye)
4. 1x1 is performing a weighted sum of the channels, so it may decide not to pick a particular feature that defines the background and not a part of the object. This is helpful as this acts like filtering. Imagine the last few layers seeing only the dog, instead of the dog sitting on the sofa, the background walls, painting on the wall, shoes on the floor, and so on. If the network can filter out unnecessary pixels, later layers can focus on describing our classes more, instead of defining the whole image.

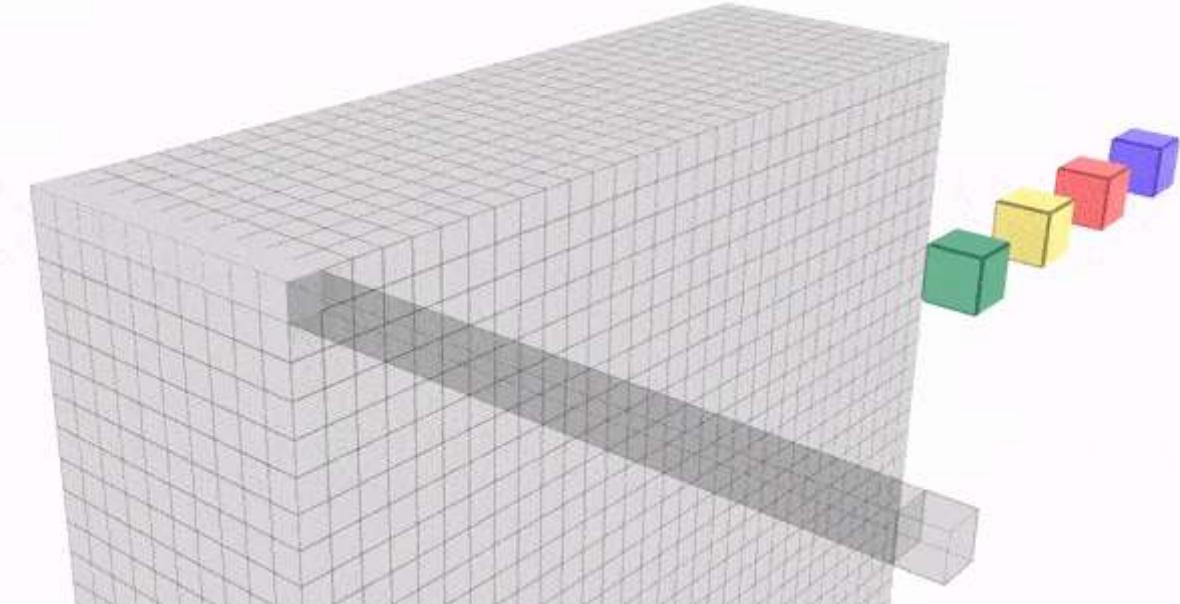
1x1 Convolutions

This is what 1x1 convolutions look like:

What you are seeing above is an input image of size $7 \times 7 \times D$ where D is the number of channels. We remember from the last lecture, that any kernel which wants to convolve, will need to possess the same number of channels. Our 1x1 kernel must have D channels. Simply put, our new kernel has D values, all defined randomly, to begin with. In 1x1 convolution, each channel in the input image would be multiplied with 1 value in the respective channel in 1x1, and then these weighted values would be summed together to create our output.



Animated 1x1:



What you see above is an input of size $32 \times 32 \times 10$. We are using 4 1×1 kernels here.

Since we have 10 channels in input, our 1×1 kernel also has 10 channels.

$$32 \times 32 \times 10 \mid 1 \times 1 \times 10 \times 4 \mid 32 \times 32 \times 4$$

We have reduced the number of channels from 10 to 4. Similarly, we will use 1×1 in our network to reduce the number of channels from 512 to 32.

Let's look at the **new** network:

$$400 \times 400 \times 3 \rightarrow (3 \times 3 \times 3) \times 32 \rightarrow 398 \times 398 \times 32 \quad \text{RF of } 3 \times 3$$

CONVOLUTION BLOCK 1 BEGINS

$$398 \times 398 \times 32 \rightarrow (3 \times 3 \times 32) \times 64 \rightarrow 396 \times 396 \times 64 \quad \text{RF of } 5 \times 5$$

$$396 \times 396 \times 64 \rightarrow (3 \times 3 \times 64) \times 128 \rightarrow 394 \times 394 \times 128 \quad \text{RF of } 7 \times 7$$

$$394 \times 394 \times 128 \rightarrow (3 \times 3 \times 128) \times 256 \rightarrow 392 \times 392 \times 256 \quad \text{RF of } 9 \times 9$$

$$392 \times 392 \times 256 \rightarrow (3 \times 3 \times 256) \times 512 \rightarrow 390 \times 390 \times 512 \quad \text{RF of } 11 \times 11$$

CONVOLUTION BLOCK 1 ENDS

TRANSITION BLOCK 1 BEGINS

MAXPOOLING(2x2)

195x195x512 ► (1x1x512)x32 ► 195x195x32 **RF of 22x22**

TRANSITION BLOCK 1 ENDS

CONVOLUTION BLOCK 2 BEGINS

195x195x32 ► (3x3x32)x64 ► 193x193x64 RF of 24x24

193x193x64 ► (3x3x64)x128 ► 191x191x128 RF of 26x26

191x191x128 ► (3x3x128)x256 ► 189x189x256 RF of 28x28

189x189x256 ► (3x3x256)x512 ► 187x187x512 RF of 30x30

CONVOLUTION BLOCK 2 ENDS

TRANSITION BLOCK 2 BEGINS

MAXPOOLING(2x2)

93x93x512 ► (1x1x512)x32 ► 93x93x32 RF of 60x60

TRANSITION BLOCK 2 ENDS

CONVOLUTION BLOCK 3 BEGINS

93x93x32 ► (3x3x32)x64 ► 91x91x64 RF of 62x62

...

RECEPTIVE FIELDS MENTIONED ABOVE ARE NOT CORRECT AND WRITTEN TO GET SOME IDEA ON HOW WAY MAY WANT THEM TO INCREASE.

Notice, that we have kept the first convolution outside of our convolution block, as now we can create a functional block receiving 32 channels and then perform 4 convolutions, giving finally 512 channels, which can then be fed to the transition block (hoping to receive 512 channels), which finally reduces channels to 32.

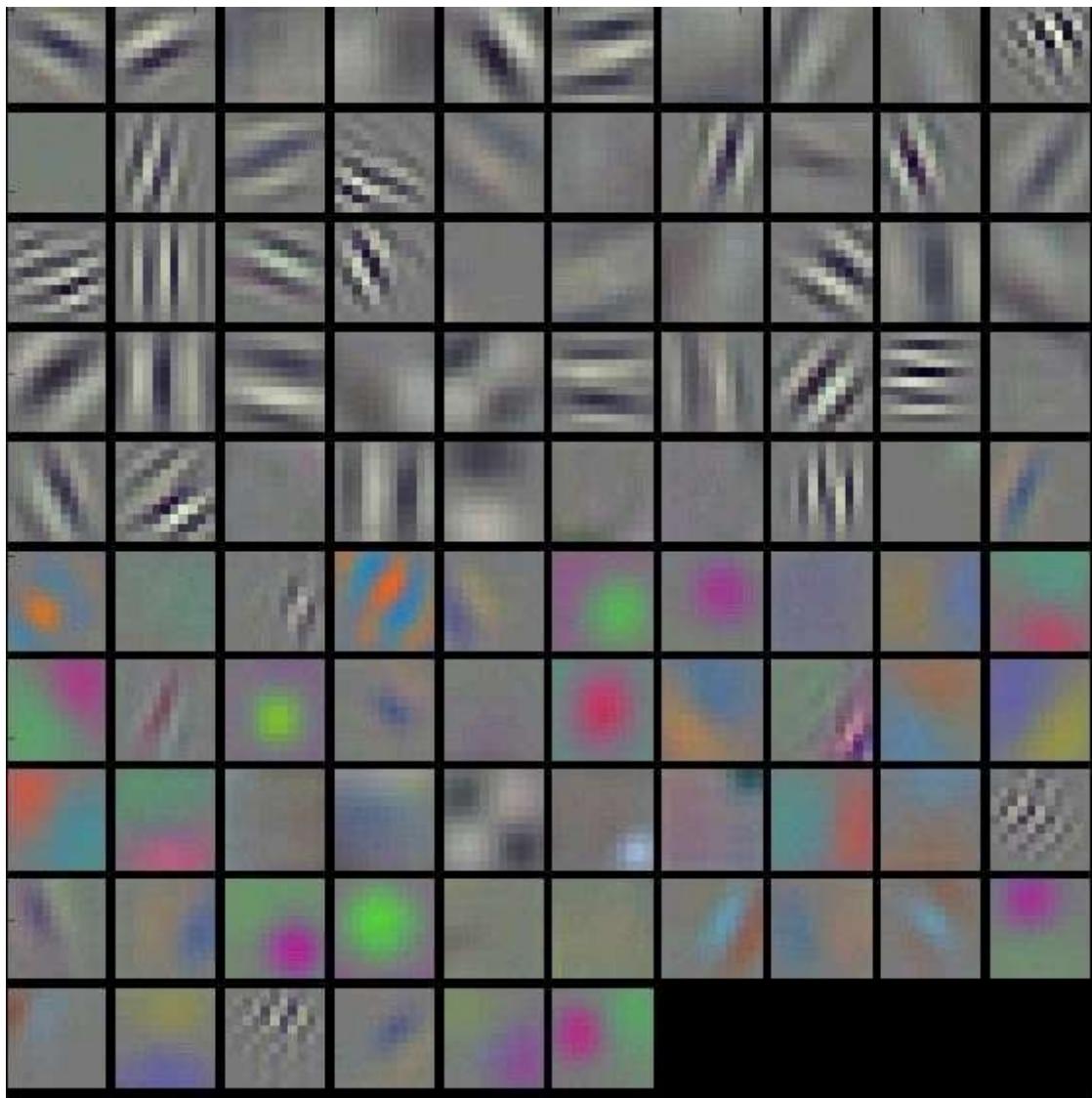
[Let's practice all of these concepts right now!](#) ↗

(<https://colab.research.google.com/drive/1oVt7T6tb90Y1EXvFaIWgm72emqZZPXi4?usp=sharing>)

Embeddings

What are embeddings?

Let's look at this image:



We know what this is, and we know how we get these visualizations by training the model.

What if we could store them and reuse them later?

For images we don't do it at this level (block 1), but certainly at block 3/4 (transfer learning). But for images, the process isn't still standardized. For NLP, however, it is!

Let's look at the same image again but with words overlapped to understand contextually what we mean.

Here I am asking you to imagine that we **can** come up with the exact representation value for each word. But something like this is tricky for the images (as there are *types* of images like natural images, astronomical, microscopic, X-ray/medical, etc). (I would also argue, that theoretically it is possible to do this for images, but for some reason, this hasn't been done)

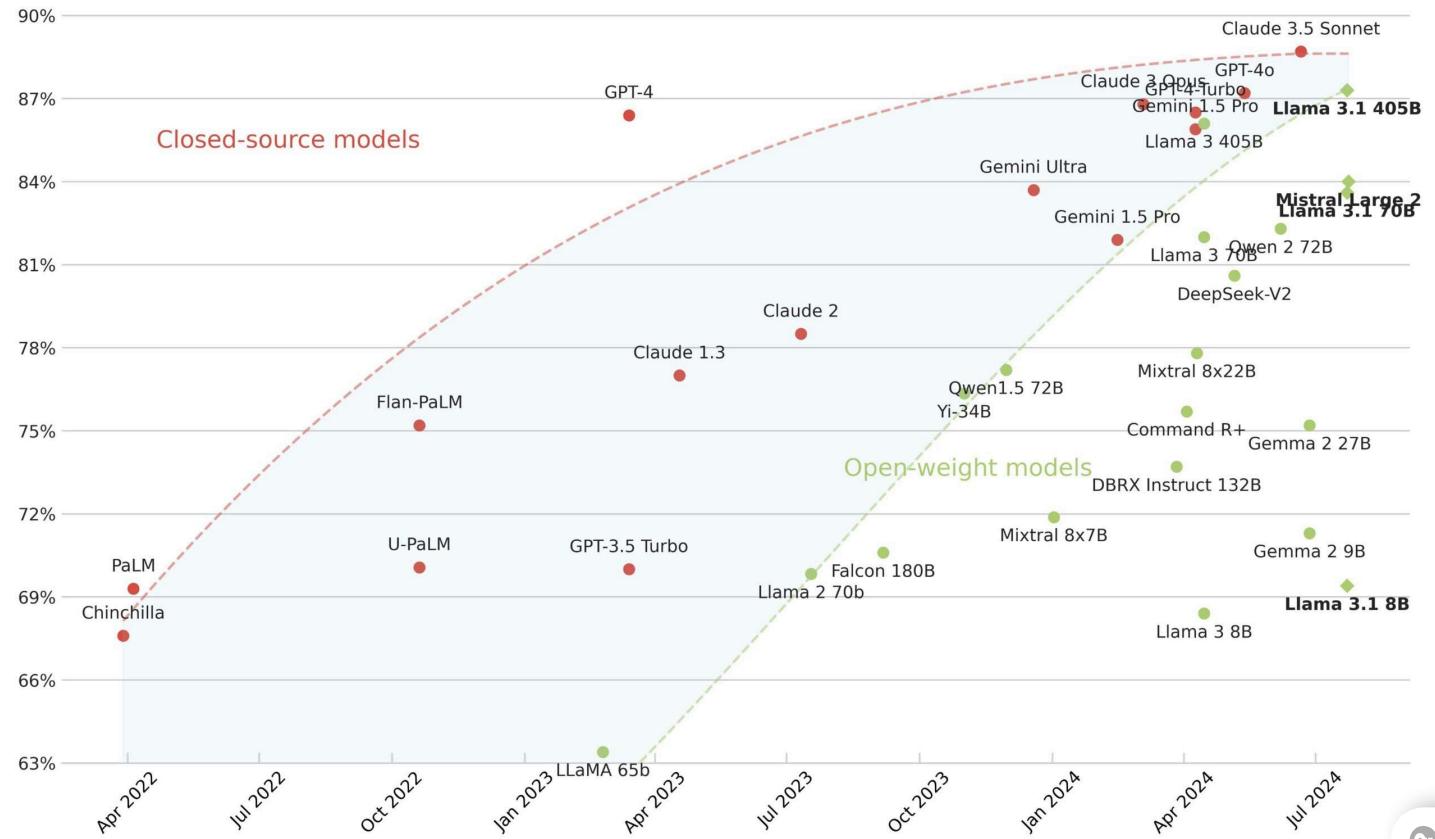
But words, well, we can get a big enough dictionary (called tokens) where we can represent every word in the world.

Closed-source vs. open-weight models

@maximelabonne

Llama 3.1 405B closes the gap with closed-source models for the first time in history.

MMLU (5-shot)



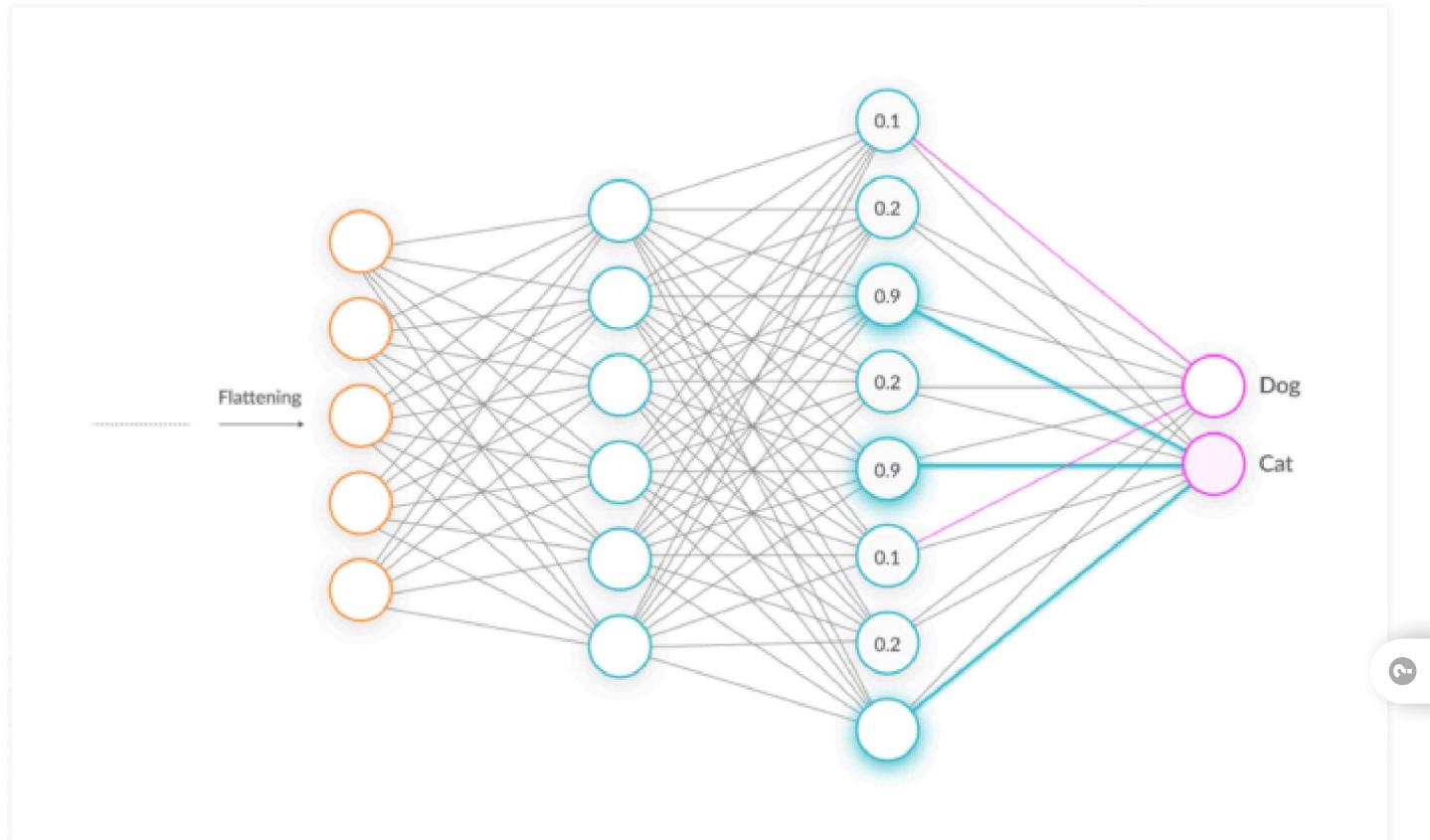
Also since we represent NLP data in 1D, it would look like this:

Positional Embeddings

What are positional embeddings?

We learned earlier that when we convert 2D data into 1D, we lose position data.

Well when we use 1D data with FC Layer, we lose position data there as well!



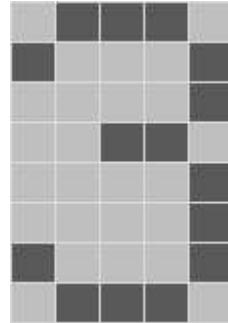
Remember: "These circles are "temporary" values that will be stored. Once you train the model, *lines are what all matters*"

But since it is easier to understand with images, let's stick to images. What digit do you think this pattern represents?



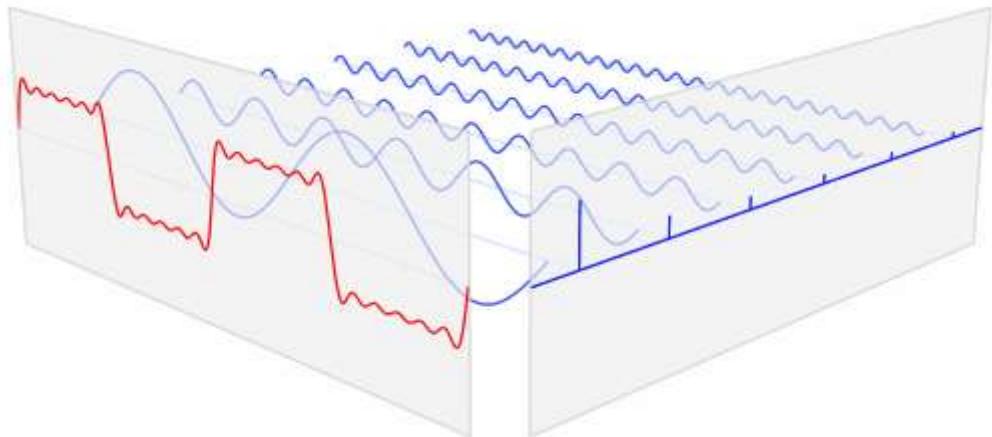
Exactly, that's the point!

This is actually a:

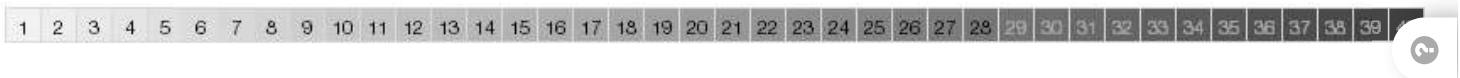


The process of conversion looks something like this, as you already know:

And as you can see we lost positional information here. How do we add the positional information back?



DNNs are surprisingly amazing at extracting the added signals to original data, just like Fourier transformers. This means, that we can add data to our original image/source and DNNs will be able to extract that information back. The convolution integral is, in fact, directly related to the Fourier transform and relies on a mathematical property of it.



buy **ADDING** the positional-related data back to each element.

But you can already see, adding 40 to one element and 1 to another element would screw up their scales (and neither this is something that FT can extract). (This is a topic for the future).

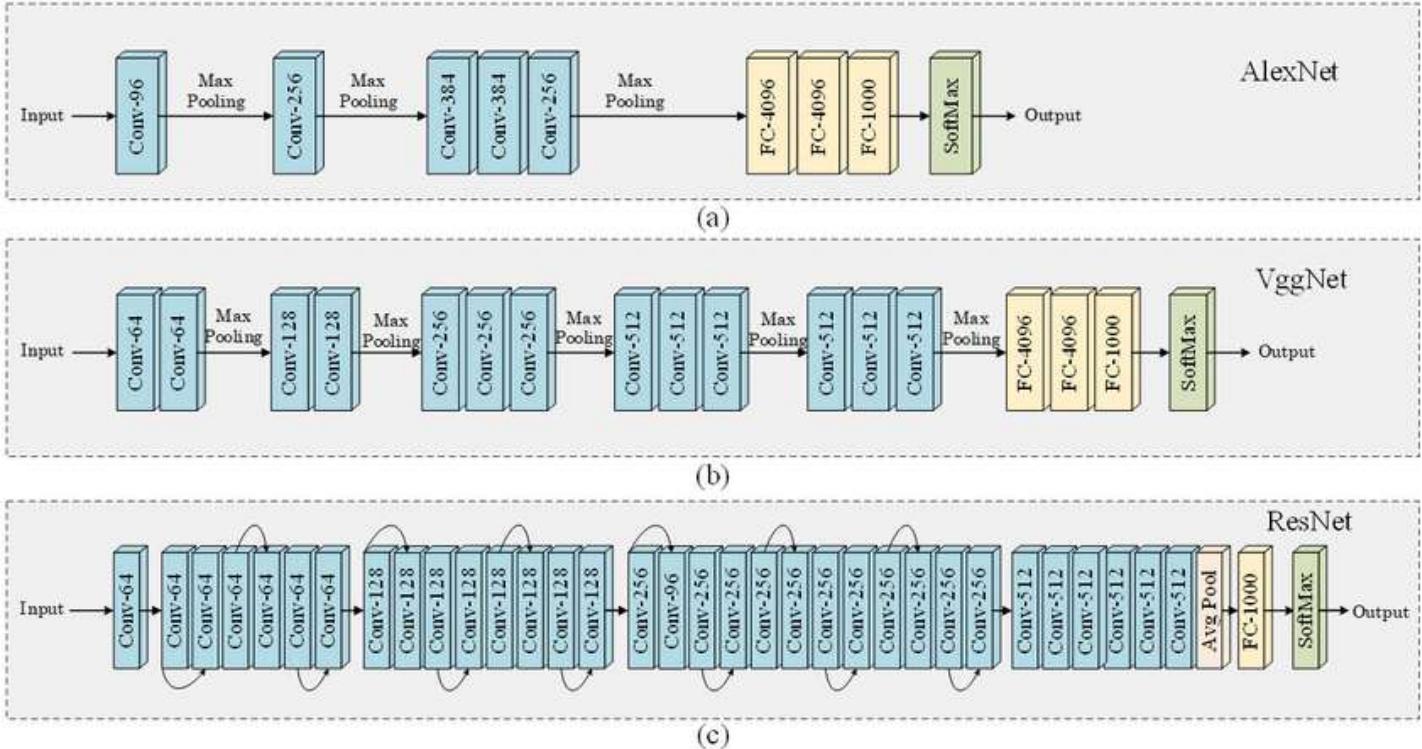
Like a word can have an embedding (say each word is represented by a 512D vector), our positions can also have embeddings (512D vector instead of just 1 number). In future sessions, it would look something like this:

The Modern Architecture

Squeeze and Expand vs Pyramidal

What we have covered as architecture is for understanding only. This architecture is called Squeeze and Expand.

Most modern architecture, however, follows Pyramidal architecture:

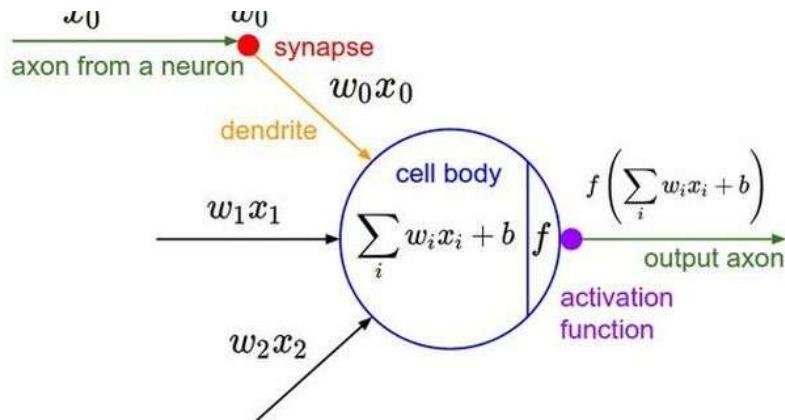


Major Points:

1. ResNet is the latest among the above. You can clearly note 4 major blocks.
2. The total number of kernels increases from $64 > 128 > 256 > 512$ as we proceed from the first block to the last (unlike what we discussed where at each block we expand to 512). Both architectures are correct, but $64 \dots 512$ would lead to lesser computation and parameters.
3. Only the most advanced networks have ditched the FC layer for the GAP layer. In TSAI we will only use GAP Layers.
4. Every network starts from 56×56 resolution!

Activation Functions and which one to use.

Their main purpose is to convert an input signal of a node in an ANN to an output signal.



For ages, this activation function has kept the AIML community occupied in the wrong direction.

After the convolution is done, we need to take a call about what to do with the values our kernel is providing us. Should we let all pass, or should we change them? This question of choice (of what to do with output) has kept us guessing. And as humans always feel, deciding what to pick and what to remove must have a complicated solution.

Why do we need an activation function?

Here is what ChatGPT thinks:

1. Activation functions allow neurons to respond in a non-linear way to the input data, which is important for modeling many real-world problems.
2. Activation functions introduce non-linearity into the network, which improves its ability to learn complex patterns and relationships in the data.

- Without activation functions, a neural network would be limited to learning only linear relationships, which would make it much less effective at solving many real-world problems.
- Activation functions are an essential component of neural networks because they enable the network to learn complex patterns and relationships in the data, and they allow the network to be more powerful and versatile.
- Activation functions are necessary for the proper functioning of neural networks, and they play a crucial role in determining the performance of the network on real-world tasks.

Let **us** understand why we need an activation function.

Let's start by asking, what is that σ doing there?

A linear activation function (or none) has two major problems:

It is not possible to use backpropagation (gradient descent) to train the model—the derivative of the function is a constant and has no relation to the input, X. So it's not possible to go back and understand which weights in the input neurons can provide a better prediction

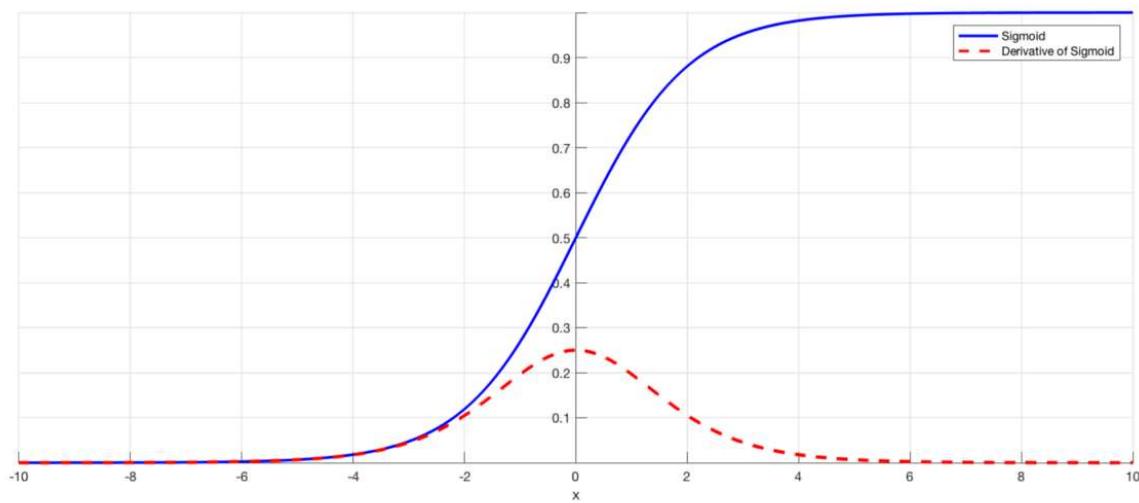
All layers of the neural network collapse into one—with linear activation functions, no matter how many layers are in the neural network, the last layer will be a linear function of the first layer (because a linear combination of linear functions is still a linear function). So a linear activation function turns the neural network into just one layer

That is why we need activation functions (other than just linear ones)

Activation Function Types

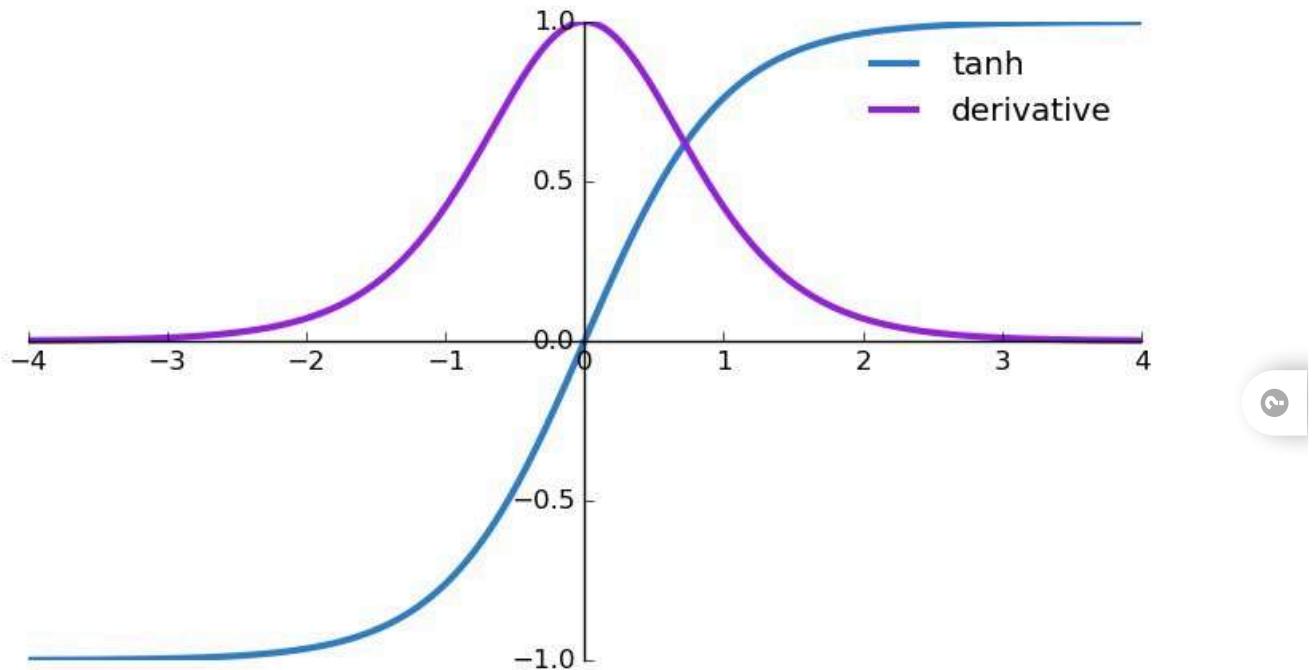
Sigmoid

We used the Sigmoid function for decades and faced something called gradient descent or gradient explosion problems. This is how Sigmoid works:



tanh

When sigmoids didn't work, we felt maybe another complicated function, called **TanH** might work:



ReLU

This is a very simple function. It allows all the positive numbers to pass on, as it is, and converts all the negatives to zero. It is a simple message to backpropagation or kernels: *"If you want some data to be passed on to the next layers, please make sure the values are positive. Negative values would be filtered out."* This also means that if some value should not be passed on to the next layers, just convert them to negatives.

ReLU Layer filters out negative numbers:

ReLU Layer

Filter 1 Feature Map

9	3	5	-8
-6	2	-3	1
1	3	4	1
3	-4	5	1



9	3	5	0
0	2	0	1
1	3	4	1
3	0	5	1

ReLU since has worked wonders for us. It is fast, simple, and efficient.

ReLU- Rectified Linear Units: It has become very popular in the past 7-8 years. It was recently (2018) proved that it had *6 times improvement in convergence* from the Tanh function. Mathematically it is just:

$$R(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

ReLU is linear (identity) for all positive values, and zero for all negative values. This means that:

It's cheap to compute as there is no complicated math. The model can, therefore, take less time to train or run.

It converges faster. Linearity means that the slope doesn't plateau, or "saturate," when x gets large.

It doesn't have the vanishing gradient problem suffered by other activation functions like sigmoid or tanh.

It's sparsely activated. Since ReLU is zero for all negative inputs, it's likely for any given unit to not activate at all.

But aren't humans known to complicate things?

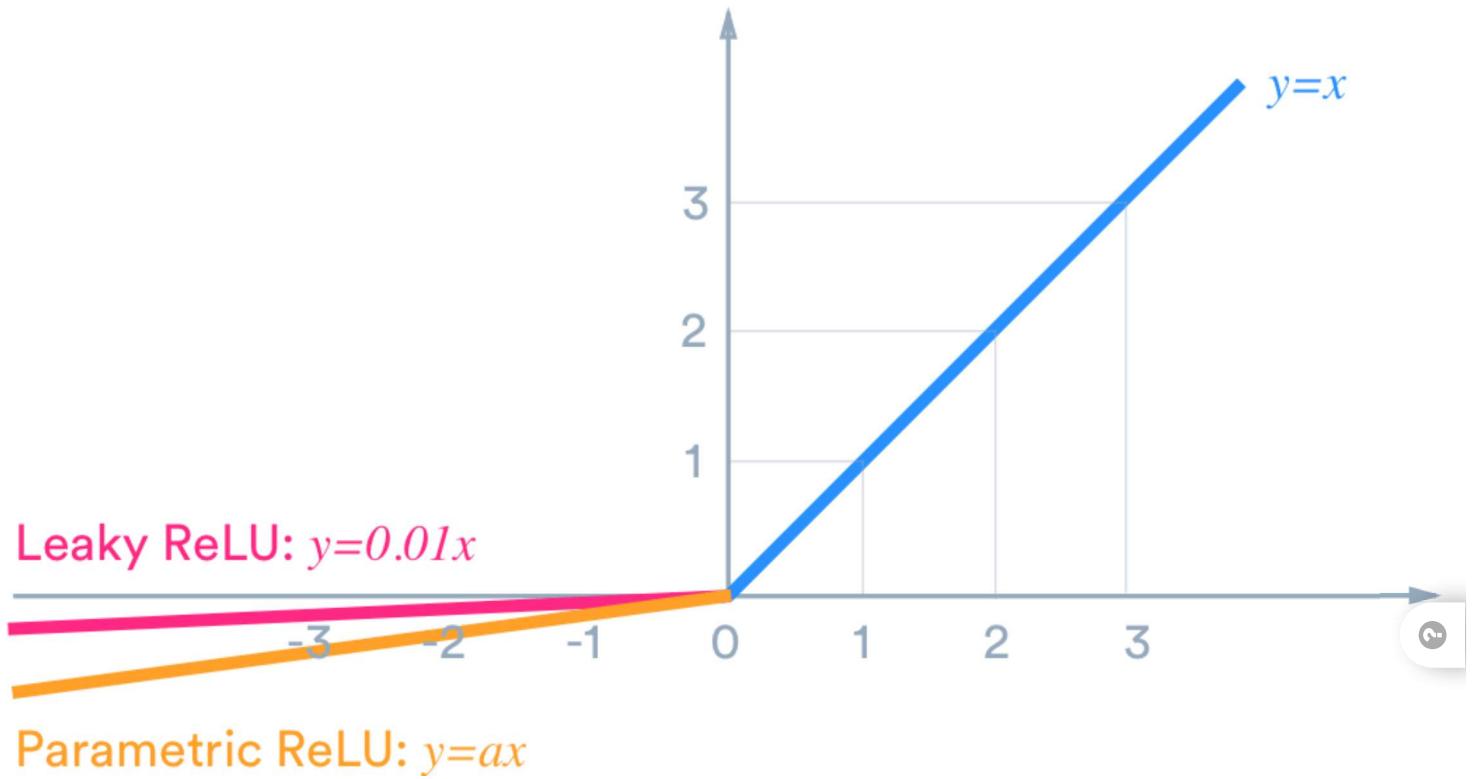
One might wonder, why such partiality to all the negative numbers. What if we allow a negative number to *leak* a few of their values?

Thence came

LeakyReLU

ReLU vs LeakyReLU

And then why stop there? You can complicate this further. Now, instead of being a Constance, we can get DNN to figure out what α should be.

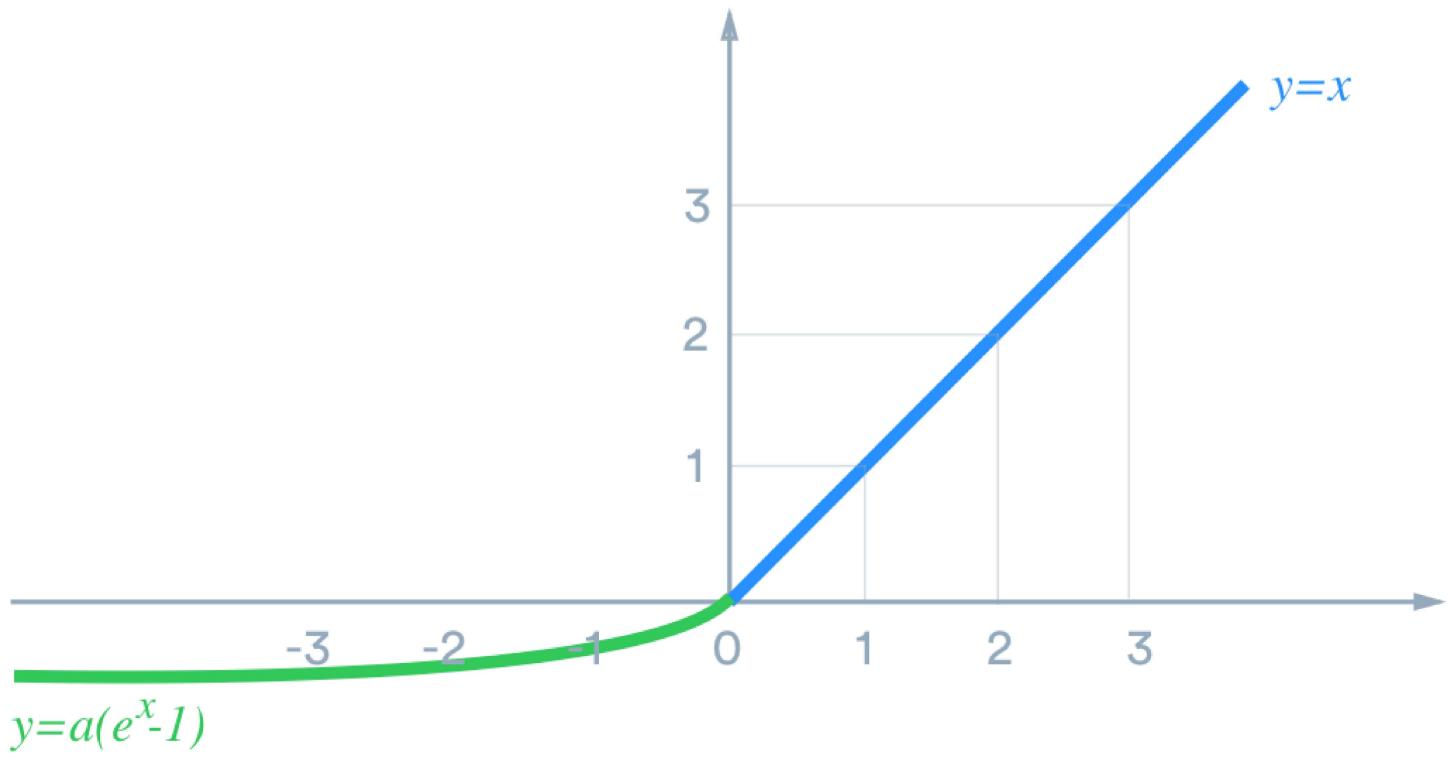


And then came a sequence of papers, each just adding a character before ELU and creating new activation functions.

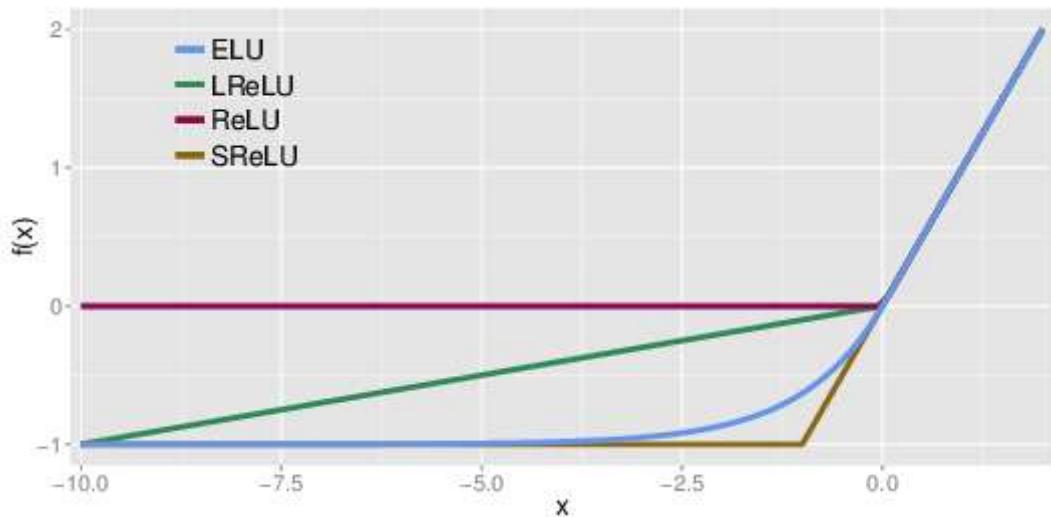
SELU



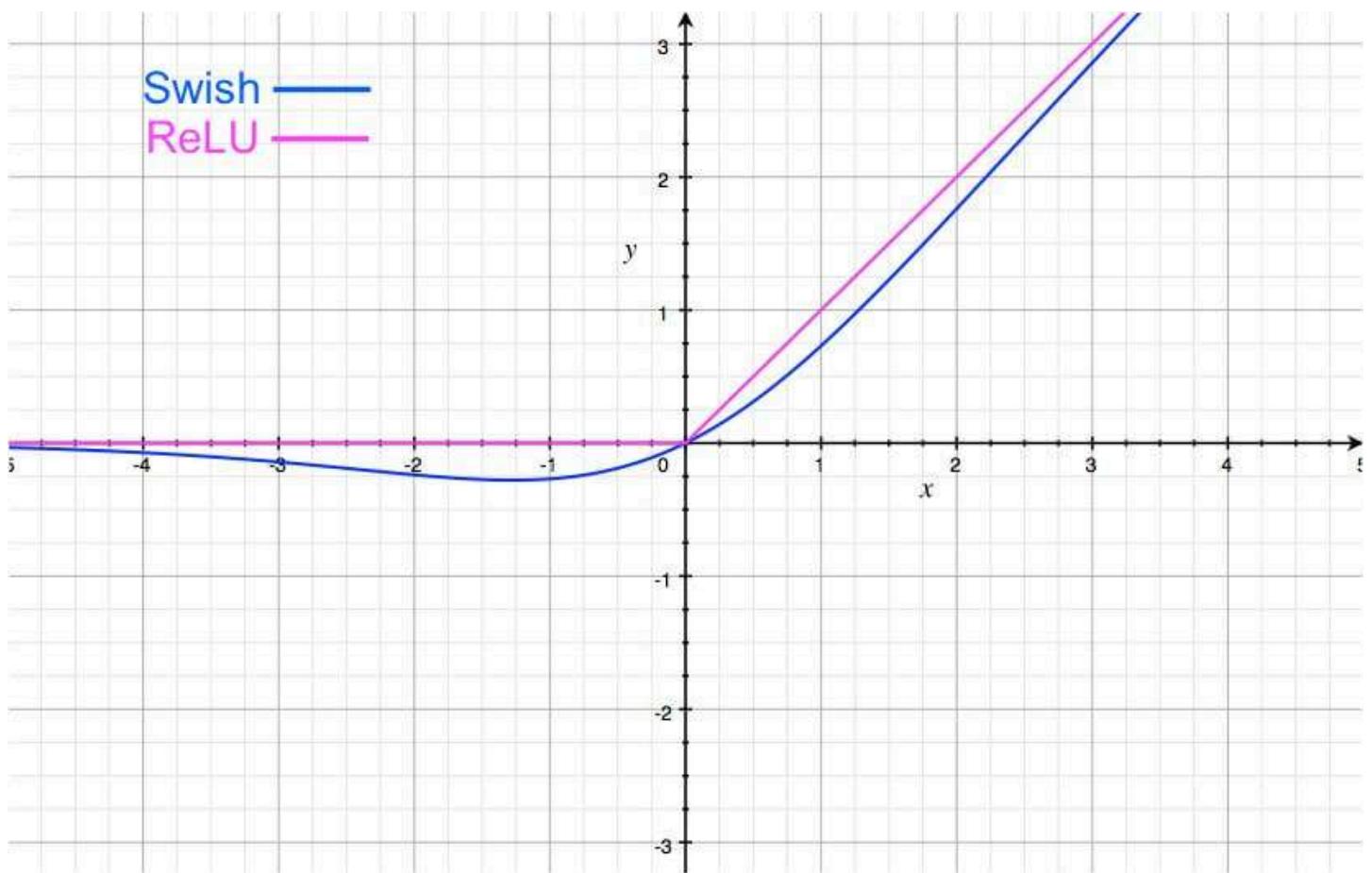
ELU



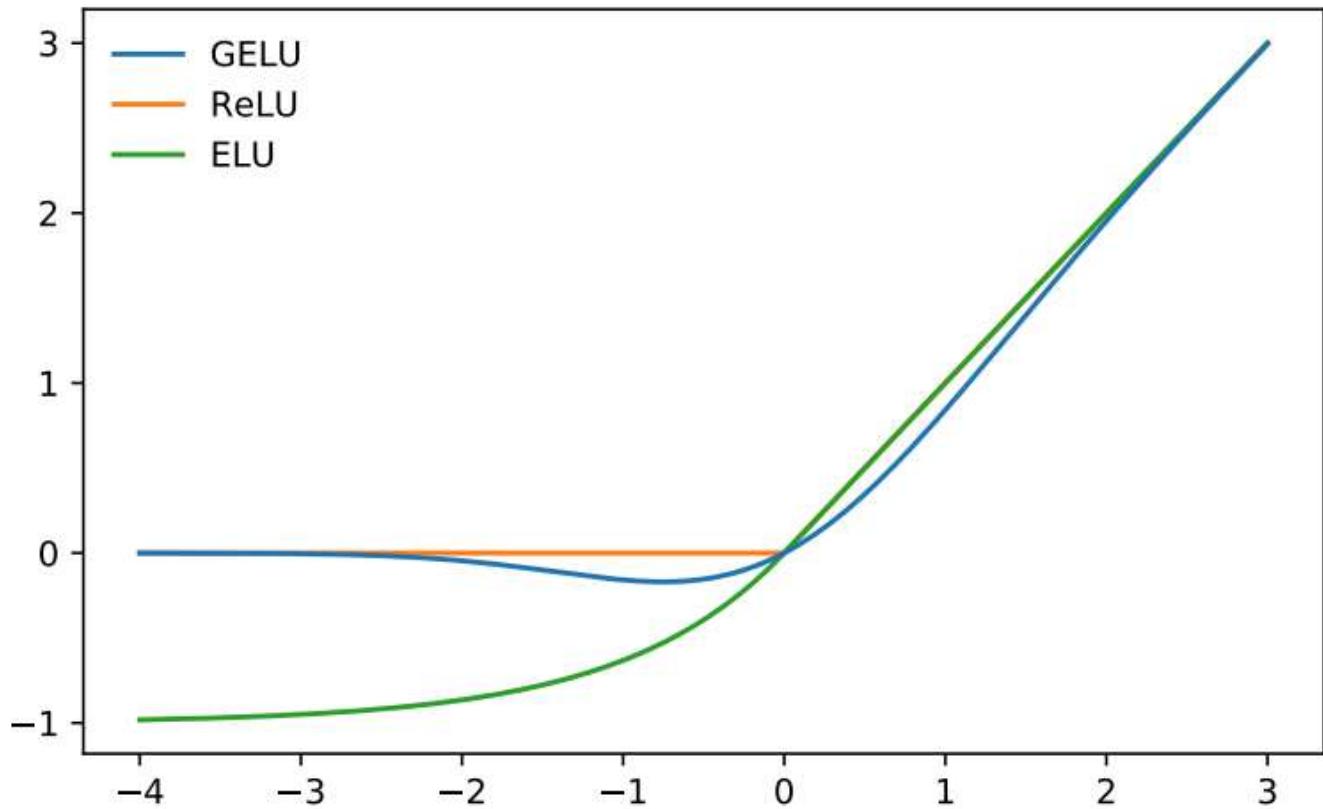
SReLU



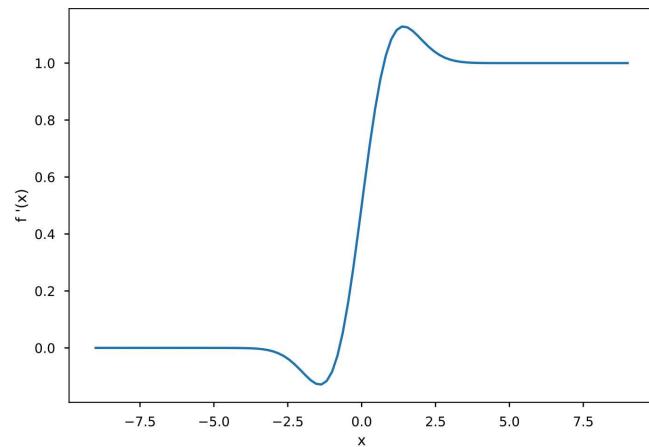
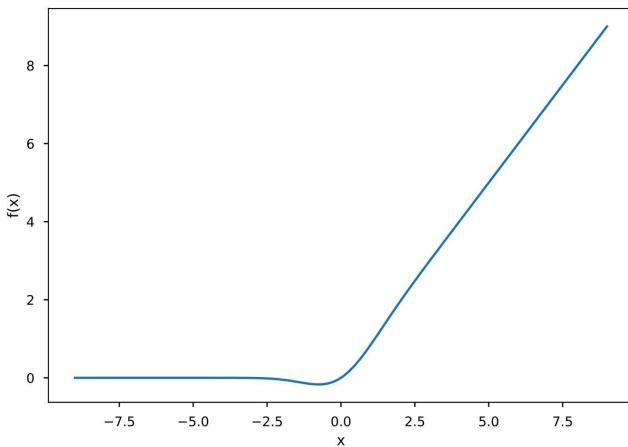
Swish



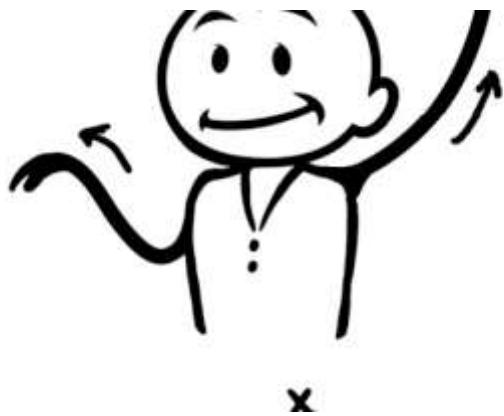
GELU



GELU function and it's Derivative



So which activation function to use?



Want to keep things simple? **ReLU**



Most of the time the innovators of these new Activation functions are using ReLU in later papers, that encouraging. There is no clear proof of which one is better. The innovators claim theirs is better, but then later peer group would release their studies that claim otherwise.

ReLU is simple, efficient, and fast, and even if any one of the above is better, we are talking about a marginal benefit, with an increase in computation. We'll stick to ReLU in our course. Also, NVIDIA has acceleration for ReLU activation, so that helps too.

Want amazing things at a cost? **GELU**

Test Error Comparison

Dataset	GELU	ReLU	ELU
CIFAR-10 (vision)	7.89%	8.16%	8.41%
CIFAR-100 (vision)	20.74%	21.77%	22.98%
TIMIT (audio)	29.3%	29.5%	29.66%
Twitter POS tagging (nlp)	12.57%	12.67%	12.91%

Median Test Error Rate after 5 training runs.

Read more [here ↗ \(https://towardsai.net/p/l/is-gelu-the-relu-successor\)](https://towardsai.net/p/l/is-gelu-the-relu-successor)



CoreSets and Curriculum Learning

CoreSets

Short for (Core Subjects) are the **minimal essence of data**. CoreSets are carefully selected, small subsets of data that preserve the performance of training on a much

larger dataset. The goal is to identify which examples actually matter for learning. TSAI IS A CORESET TEACHING :)

Why CoreSets matter:

- Training on huge datasets is computationally expensive and includes lots of redundant or low-value data
- CoreSets aim to reduce training time and cost without sacrificing performance
- They can also de-bias training by filtering out noisy or overrepresented patterns

CoreSets come from optimization theory and active learning. Imagine training a model with 1 billion examples, but 90% of the model's performance could be achieved with just 1 million well-chosen examples!

They use ideas like:

- Diversity Samples: Choose examples that are representative but not repetitive
- Gradient-based selection: Pick samples that maximize change in model weights
- Loss-based filtering: Keep only the examples that challenge the model the most.

Curriculum Learning

Teaching the models Like We Teach Kids.

CL takes inspiration from the human education system. Start simple, then go complex.

Core Idea: Instead of feeding the model random data in arbitrary order, curriculum learning structures the training process. The model first learns from easy examples, then progressively tackles harder and more nuanced ones.

Benefits of CL:

- Faster Convergence: models learn faster when they aren't overwhelmed
- Better Generalization: learned features build on top of simple ones
- Higher stability: in early training (fewer spikes in loss or divergence)

What Makes a Curriculum?

- **Perplexity:** Start with data that the model finds easy to predict
- **Length:** Begin with short and structured examples (text)
- **Complexity:** Gradually introduce irony, multi-hop reasoning, rare vocabulary, etc
- **Domain transition:** Move from General Knowledge to specialized fields (e.g. from Encyclopedia to scientific papers)

Together: CoreSets × Curriculum Learning

- **Efficient learning** on reduced data
- **Robust model behavior** even in a complex domain
- **\$\$\$ Help save massive compute cost \$\$\$**

References

1. [Interactive Introduction to Fourier Transforms](https://www.v7labs.com/blog/neural-networks-activation-functions)
2. [Why do neural networks work so well - Quora](https://www.jezzamon.com/fourier/)
3. [What is a Convolutional Neural Network](https://poloclub.github.io/cnn-explainer/)
4. [Feature Visualization - How Neural networks build up their understanding of images](https://distill.pub/2017/feature-visualization/)
5. [Computing Receptive Fields](https://distill.pub/2019/computing-receptive-fields/)
6. [Deconvolution and Checkerboard artifacts](https://distill.pub/2016/deconv-checkerboard/)
7. [Activation functions in Neural Networks](https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6)
8. [Is GELU, the ReLU successor?](https://towardsai.net/p/l/is-gelu-the-relu-successor)
9. [Activation Functions in Deep Learning: A Comprehensive Survey and Benchmark](https://arxiv.org/pdf/2109.14545.pdf)

The Assignment

1. **Assignment 1:** Make an MNIST-based model that has the following characteristics:
 1. Has fewer than 25000 parameters
 2. Gets to **test** accuracy of 95% or more in 1 Epoch
2. Once done (locally or on Google Colab), upload to GitHub, and make a README file explaining your architecture and Logs.
3. Once you achieve the above results,
 1. Share the link to your README.md file, it must be a beautifully formatted file.
[250]
4. **Assignment 2:**
 1. Repeat my Class Exercise, and then take screenshot of your Scheduler Logs, and show Gemini output atleast twice in same photo [250]
 2. Move your last S3 Submission to Lambda and share the web URL which would allow me to trigger the usage/your-app (no EC2!!)! [500][Optional]

Videos

Studio ([Transcript \(\[https://canvas.instructure.com/courses/12597696/files/310808465?\]\(https://canvas.instructure.com/courses/12597696/files/310808465?wrap=1\)](https://canvas.instructure.com/courses/12597696/files/310808465?wrap=1))

[Download \(\[https://canvas.instructure.com/courses/12597696/files/310808465/download?download_frd=1\]\(https://canvas.instructure.com/courses/12597696/files/310808465/download?download_frd=1\)\)](https://canvas.instructure.com/courses/12597696/files/310808465/download?download_frd=1)

ERA V4 Session 4 Studio



GMeet

ERA V4 Session - 4 GMeet



