# ELV26 Tutorial 3: **Video Learning**

Ellis Brown
2026-02-10

# Overview
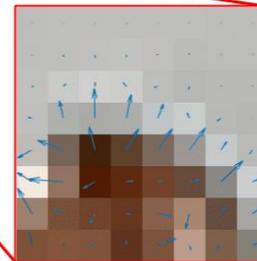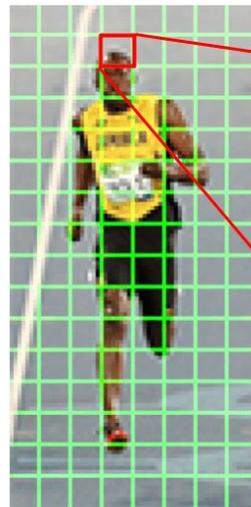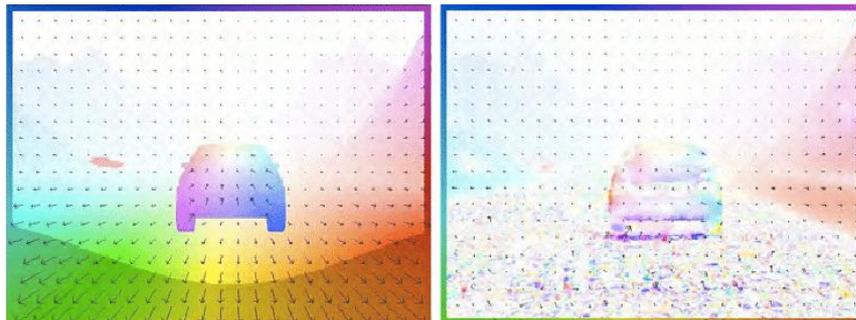
I. Context
II. Working w/ Video Data
III. Towards Video-LLMs
IV. Practical Training
V. Video-LLM failure modes

# I. Context

# Birds-eye-view



**Era 1 — Hand-Crafted Motion**

- **Optical Flow**, hand-crafted features (eg, Dense Trajectories)
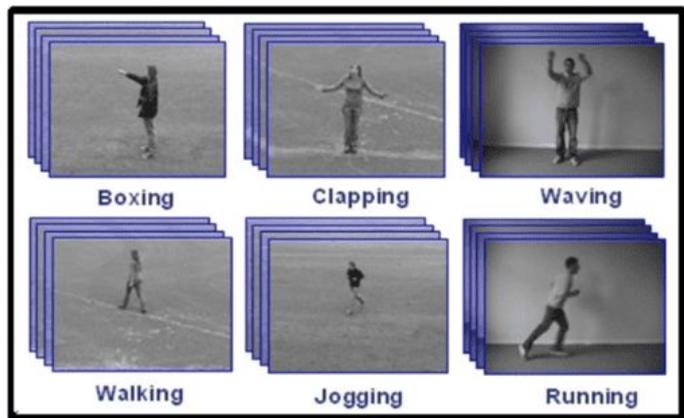- Track pixel movements
- ❌ low-level, no semantic understanding

**~2000**

**~2014**

# Birds-eye-view



Kinetics-400

Juggling soccer ball

Pushing car

Opening something

Something-Something V1&V2

Squeezing something



Boxing  Clapping  Waving

Walking  Jogging  Running

~2014

~2022

**Era 2 — Deep Classification**

- **3D** CNNs, Video-SSL
- Task: N-way classification (eg Kinetics-400 human action recognition)
- → shortcut: classify w/ one frame…

*Video-MME*

What is the woman wearing when standing in front of the mirror?

A. A green suit with a purple blouse.          B. A sleeveless, bright pink dress.

C. A bright blue blouse with a long, red cape.          D. A bright yellow coat.

[Option C]          [Option A]          [Option B]          [Option D]

NENDEZ-VOUS

01:10          04:12          27:52          31:16

(eg, Dense Trajectories)

**Object Count**

How many chairs are there in this room?

*Answer: 4*

**Relative Distance**

Measuring from the closest point of each object, which of these objects (refrigerator, sofa, ceiling light, cutting board) is the closest to the printer?

*A. refrigerator  B. sofa  C. ceiling light  D. cutting board*

**Appearance Order**

What will be the first-time appearance order of the following categories in the video: basket, printer, refrigerator, kettle?

*A. kettle, basket, printer, refrigerator*
*B. refrigerator, printer, basket, kettle*
*C. basket, printer, refrigerator, kettle*
*D. basket, refrigerator, kettle, printer*

**Relative Direction**

If I am standing by the refrigerator and facing the sofa, is the kettle to my left, right, or back?

*A. left  B. right  C. back*

**Object Size**

What is the length of the longest dimension (length, width, or height) of the refrigerator in centimeters?

*Answer: 119*

**Absolute Distance**

Measuring from the closest point of each object, what is the distance between the bed and the sofa in meters?

*Answer: 3.2*

**Room Size**

What is the size of this room (in square meters)? If multiple rooms are shown, estimate the size of the combined space.

*Answer: 57.6*

**Route Plan**

You are a robot beginning at the toilet and facing the washer. Navigate to the pan. Fill in this route: 1. Go forward until the washing machine 2. [?] 3. Go forward until the sofa 4. [?] 5. Go forward until the pan.

*A. Turn Left, Turn Left  B. Turn Left, Turn Right*
*C. Turn Back, Turn Right  D. Turn Right, Turn Right*

Video Starting Position

Camera Trajectory

Video Ending Position

Camera Perspective

VSI-Bench

**Era 3 — Generative *Understanding***

- Labels → *explanations*
- *Reasoning* about time, physics, intent
- Strong semantic understanding, *poor basic physical intuition*

**~2022**

# Birds-eye-view

**Era 1 — Hand-Crafted Motion**

- **Optical Flow**, hand-crafted features (eg, Dense Trajectories)
- Track pixel movements
- ❌ low-level, no semantic understanding

**Era 3 — Generative *Understanding***

- Labels → *explanations*
- *Reasoning* about time, physics, intent
- Strong semantic understanding, *poor basic physical intuition*

~2014

**~2000**

**~2022**

**Era 2 — Deep Classification**

- **3D** CNNs, Video-SSL
- Task: N-way classification (eg Kinetics-400 human action recognition)
- → shortcut: classify w/ one frame…

# II. Working w/ Video Data

# General deep learning process

Dataloading

1. Load raw data
2. Preprocess data
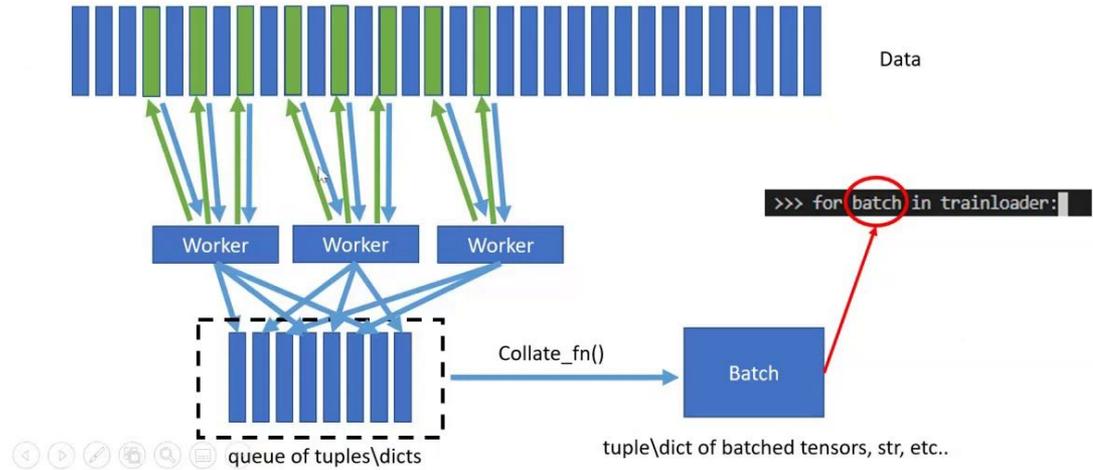3. Gather data samples into batches

Model training

1. Load data batch onto GPU
2. Perform model forward pass and loss computation
3. Perform backpropagation to update model parameters

# Dataloading

Data workers run on CPU cores

Runs asynchronously while model training is done on GPU
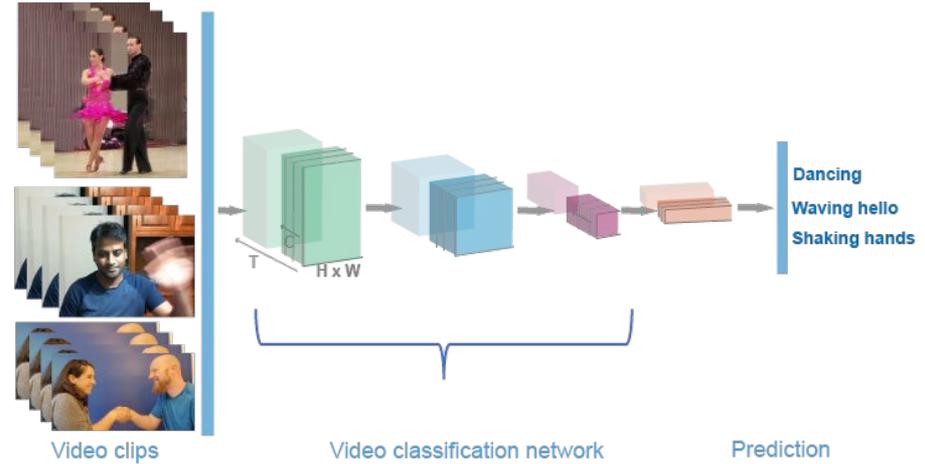
## Inside the Dataloader



Data

Worker   Worker   Worker

Collate_fn()

Batch

queue of tuples\dicts

tuple\dict of batched tensors, str, etc..

```
>>> for batch in trainloader:
```

https://www.youtube.com/watch?v=Sj-gIb0QiRM

# Videos

Sequences of image frames

Store key frames and "deltas" between key frames to achieve storage (lossy) compression

Enable learning temporal information, motion, object behavior, world models



T — C — H x W

Dancing
Waving hello
Shaking hands

Video clips     Video classification network     Prediction

# Video datasets

Kinetics400: clips of human actions

Walking Tours: walking city tours

SomethingSomething: clips of human actions

Ego4D: egocentric daily life videos

BDD100K: driving dashcams

Waymo Open: driving camera videos

# Video datasets sizes

| Dataset | Domain | Ego | Pre | Bal | Annot | Avg. Dur (sec) | Dur (hr) | #Videos | Frame Resolution |
|---|---|---|---|---|---|---|---|---|---|
| *Diverse Pretraining* | | | | | | | | | |
| Kinetics-400 (Kay et al., 2017) | Actions | ✗ | ✓ | ✓ | Class | 10.2 | 851 | 400 | 340 × 255 |
| WebVid-2M (Bain et al., 2021) | Open | ✗ | ✓ | ✗ | Weak | 18 | 13k | – | 320 × 240 |
| HowTo100M (Miech et al., 2019) | Instructions | ✗ | ✓ | ✗ | Weak | 4 | 135k | – | – |
| *Egocentric* | | | | | | | | | |
| Epic-Kitchens (Damen et al., 2022) | Cooking | ✓ | ✗ | ✗ | Loc. | 510 | 100 | 37 | 1920 × 1080 |
| Ego-4D (Grauman et al., 2022) | Daily | ✓ | ✗ | ✗ | Loc. | 1446 | 120 | 931 | 1920 × 1080 |
| Meccano (Ragusa et al., 2023) | Industry | ✓ | ✗ | ✗ | Loc. | 1247 | 849 | 20 | 1920 × 1080 |
| Assembly-101 (Sener et al., 2022) | Assembly | ✓ | ✗ | ✗ | Loc. | 426 | 167 | 362 | 1920 × 1080 |
| *ImageNet-aligned* | | | | | | | | | |
| R2V2 (Gordon et al., 2020) | ImageNet | ✗ | ✓ | ✓ | Class | – | – | – | 467 × 280 |
| VideoNet (Parthasarathy et al., 2022) | ImageNet | ✗ | ✓ | ✓ | Class | 10 | 3055 | – | |
| Walking Tours (ours) | Urban | ✓ | ✓ | ✗ | None | 5880 | 23 | 10 | 3840 × 2160 |

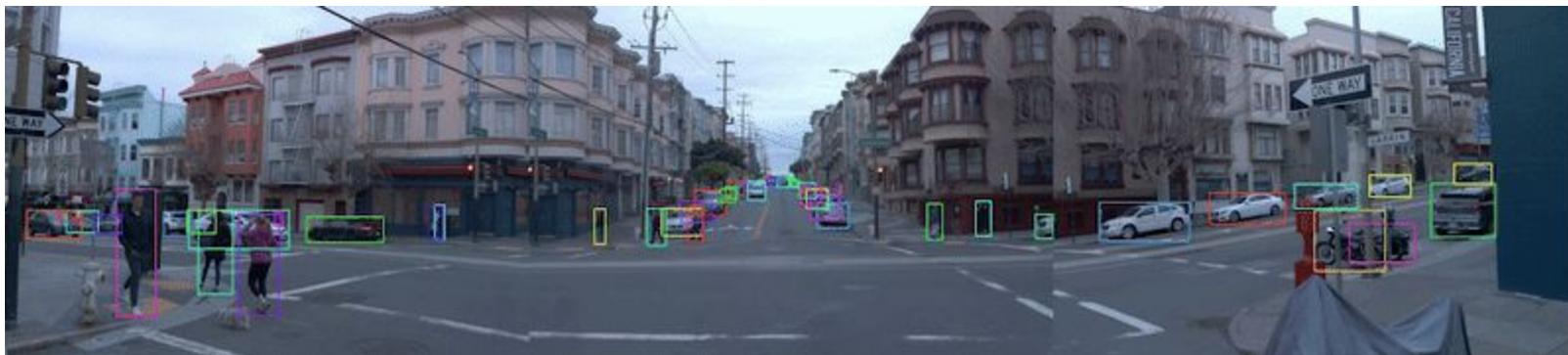For reference, ImageNet is 1.3M images of size ~470x390

# Example: WalkingTours

# EPIC-Kitchens

# Example: Waymo Open

# Video Processing

Storage - time tradeoff:

1. Videos are compressed stacks of images
2. It takes time to decode videos into data arrays

Frame sampling

1. How much time in-between frames?
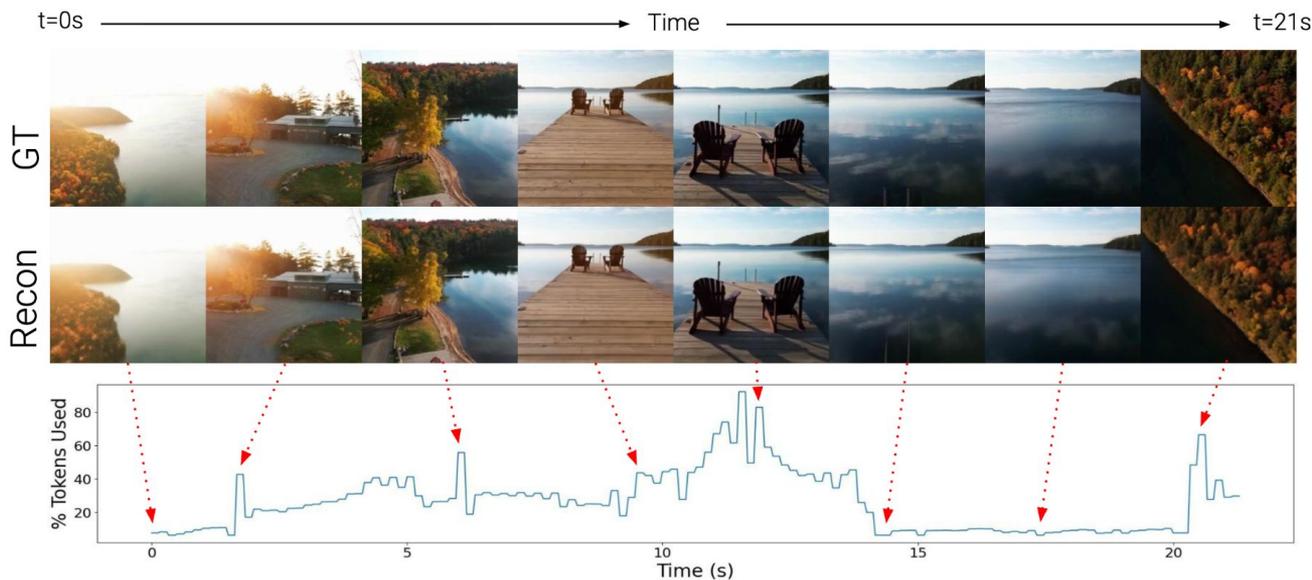2. What level of temporal granularity do you care about?

# Frame sampling



$\Delta t = 0$ (0s)　　$\Delta t = 15$ (0.5s)　　$\Delta t = 30$ (1s)　　$\Delta t = 45$ (1.5s)

# Advanced: Adaptive Tokenization



**Figure 1    ElasticTok adaptively represent video based on information available**. (Top) Ground-truth video frames. (Middle) Reconstructed frames with varying token usage. (Bottom) The bottom section depicts how ElasticTok dynamically adjusts token allocation over time, with the percentage of tokens used correlating to different content complexities in the video.
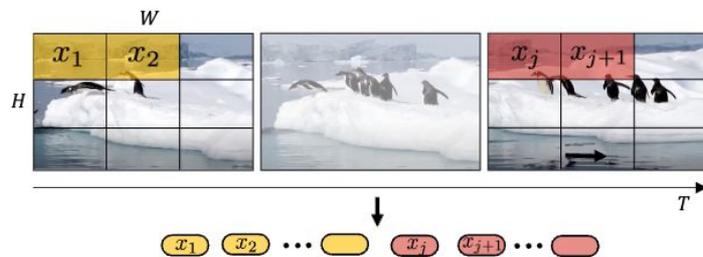
ElasticTok: Adaptive Tokenization for Image and Video

# Advanced: Video Tokenization



Figure 2: Uniform frame sampling: We simply sample $n_t$ frames, and embed each 2D frame independently following ViT [18].
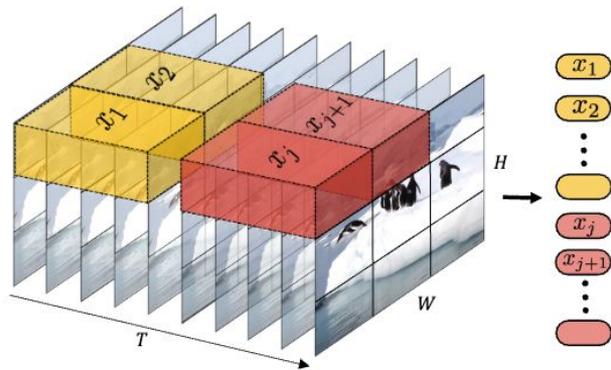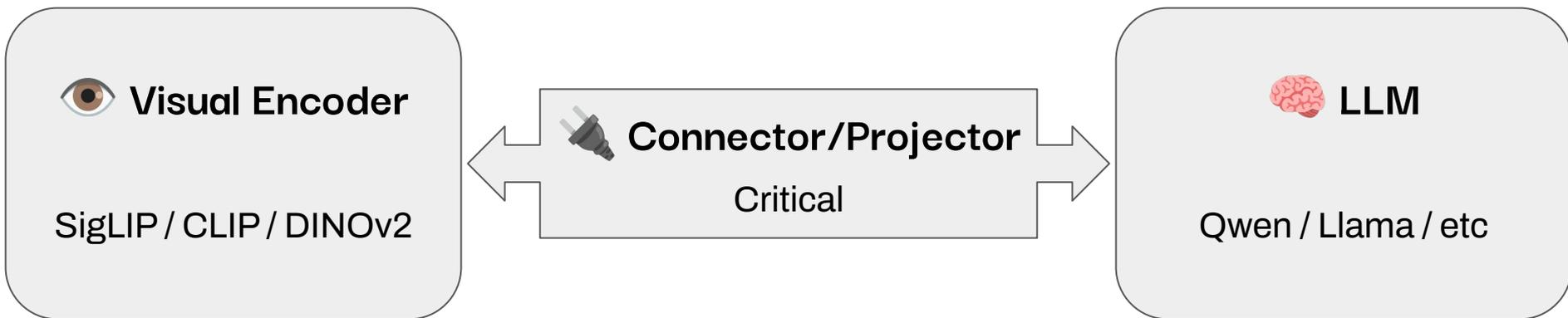

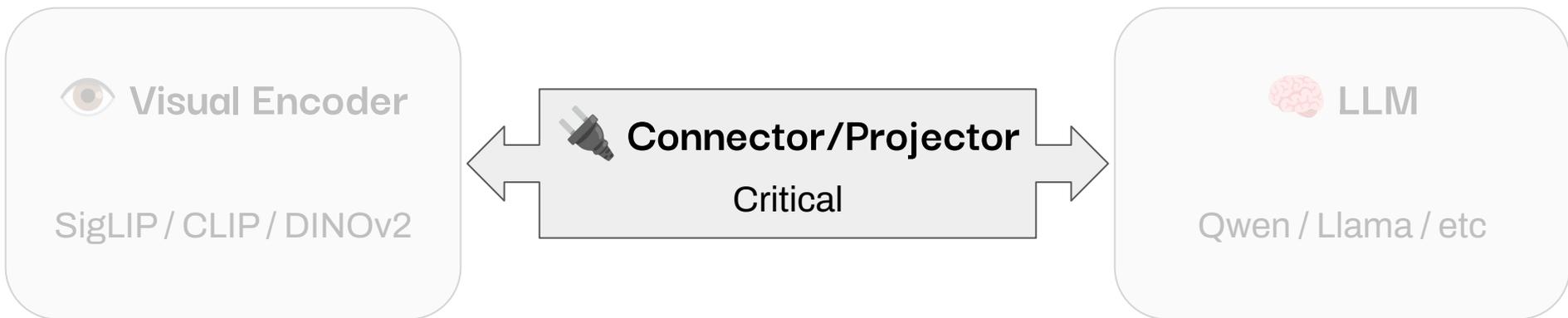
Figure 3: Tubelet embedding. We extract and linearly embed non-overlapping tubelets that span the spatio-temporal input volume.

[2103.15691] ViViT: A Video Vision Transformer

# III. Towards Video-LLMs

# High-level of MLLMs

**👁 Visual Encoder**

SigLIP / CLIP / DINOv2

**🔌 Connector/Projector**

Critical

**🧠 LLM**

Qwen / Llama / etc

# High-level of MLLMs

# Flamingo
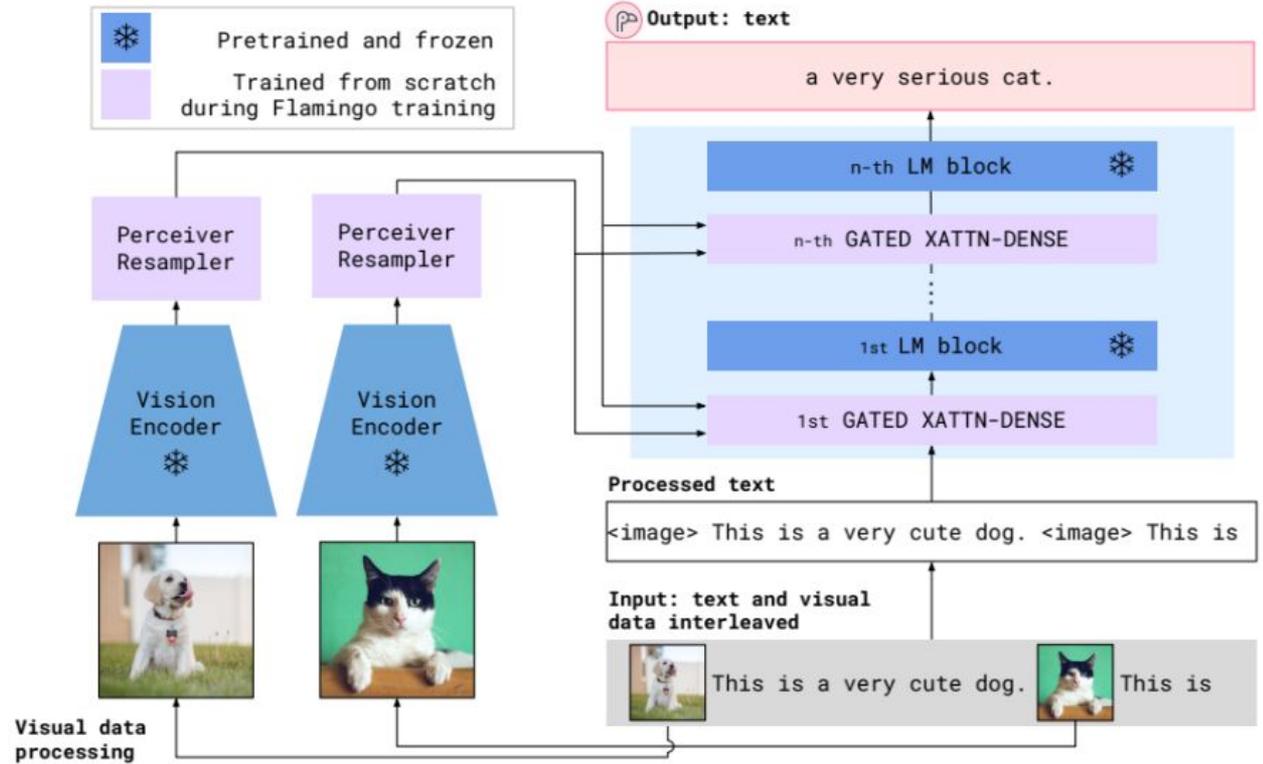
- Very flexible!
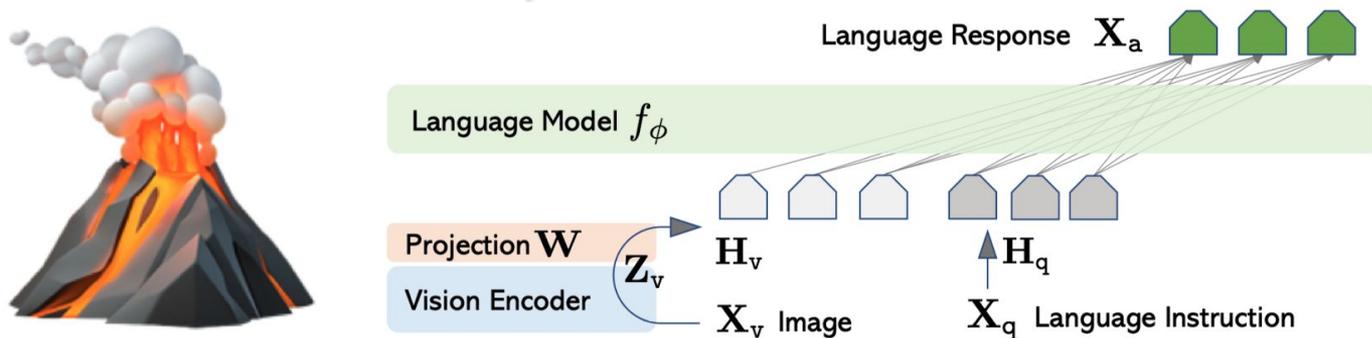- xattn to visuals
- Too heavy



Figure 3 | **Overview of the Flamingo model.** The Flamingo models are a family of visual language model (VLM) that can take as input visual data interleaved with text and can produce free-form text as output. Key to its performance are novel architectural components and pretraining strategies described in Section 3.
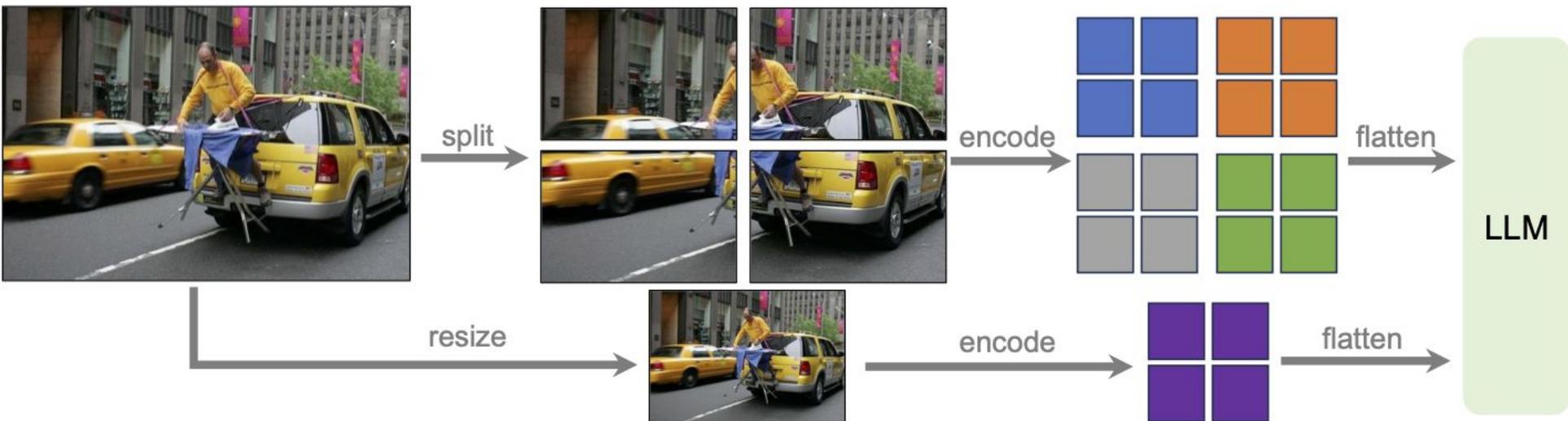
[2204.14198] Flamingo: a Visual Language Model for Few-Shot Learning
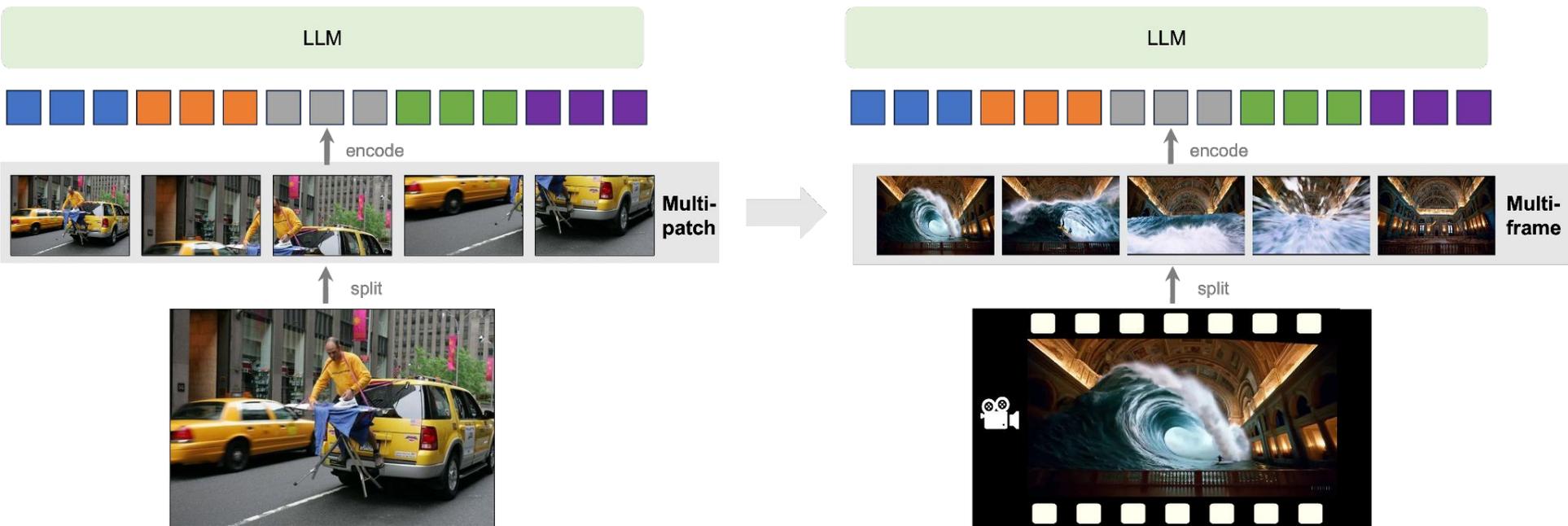
# LLaVA: Connect Single-Image to LLM

"Visual Instruction Tuning"

→ just tune the projector (linear layers). *light!*

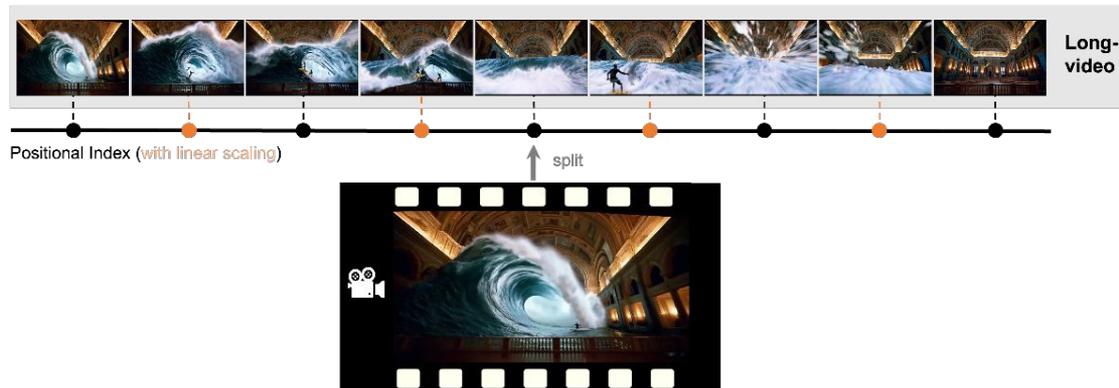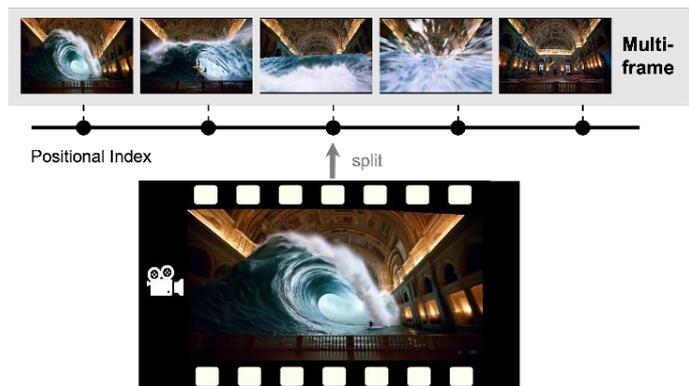# High-res processing: pass multiple crops

# Multi-Crop (1 img) → Multi-*frame* (*video*)

# Multi-Crop (1 img) → Multi-*frame* (*video*)

- Process more frames at test time
  - Miss fewer frames! But *still* drop many

# Video-LLM Landscape

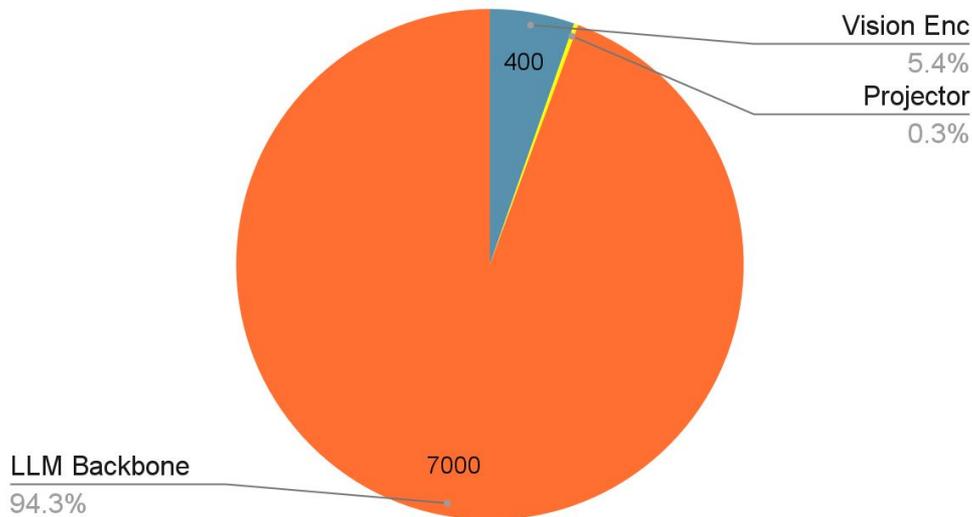| Organization | Models |
|---|---|
| **Proprietary** | |
| Google | Gemini 3.0 Pro, 2.5 Pro/Flash |
| OpenAI | GPT-5/5.2, GPT-4o |
| xAI | Grok 4/4.1 |
| Anthropic | Claude 4.5 (frames only) |
| Reka | Core, Flash, Edge |
| **Open — General** | |
| Alibaba | Qwen3-VL, Qwen2.5-VL, QVQ-72B |
| OpenGVLab | InternVL3, InternVL 3.5 |
| Meta | Llama 4 Scout/Maverick, Llama 3.2 Vision |
| Google | Gemma 3 |
| Moonshot AI | Kimi-VL, Kimi-VL-Thinking |
| OpenBMB | MiniCPM-V 4.5, MiniCPM-o 2.6 |
| Zhipu AI | GLM-4.6V |
| Microsoft | Phi-4 Multimodal |
| DeepSeek | DeepSeek-VL2, Janus-Pro |
| ByteDance | Seed1.5-VL, BAGEL |
| **Open — Video-First** | |
| LLaVA Labs | LLaVA-OneVision 1.5, LLaVA-Video |
| DAMO-NLP-SG | VideoLLaMA 3, VideoLLaMA 2 |
| OpenGVLab | VideoChat-Flash, VideoChat2 |
| ByteDance | Vidi 2.5 |
| Meta et al. | Apollo |
| AI2 | Molmo |

# IV. Practical Training

# MLLM Training ≈ LLM Training

- LLM: 7B
- Vision: 400m
- Projector: ~20M

**Good News:** We can leverage the entire ecosystem of LLM optimization (FlashAttention, vLLM, FSDP, etc.)

Parameters



Vision Enc
5.4%

Projector
0.3%

400

7000

LLM Backbone
94.3%

# Reducing Quadratic Cost of Attention:
# **Flash Attention**

Math:

- 1 Image ≈ 576 tokens.
- 1 min video (Sampled 1fps) = 60 images ≈ 34,560 tokens.
- Standard Attention cost: O(N^2)

Solution: **FlashAttention-2**

- <u>What it is:</u> *IO-aware* exact attention. Computes attention by tiling blocks to avoid HBM (High Bandwidth Memory) reads/writes.
- <u>Impact:</u> Linear-scaling memory usage.
- <u>Practical Advice:</u> Try loading HF transformer model with **attn_implementation="flash_attention_2"** →→→

```python
import torch
from transformers import AutoModelForCausalLM

model = AutoModelForCausalLM.from_pretrained(
    "your-model-id",
    torch_dtype=torch.bfloat16
    attn_implementation="flash_attention_2"
).to("cuda")
```

# Handling Large Batches:
# **Gradient Accumulation**

**Problem:**
- Video-LLMs are VRAM-heavy. A single GPU often fits only `batch_size=1` or `2`.
- Small batches cause noisy gradients and unstable training.

**Solution:**
- **Virtual Batch Size:** Decouple your *hardware limit* from your *training batch size*.
- **Mechanism:**
  1. Forward/Backward pass on Micro-Batch 1
     → **Do not clear grads**
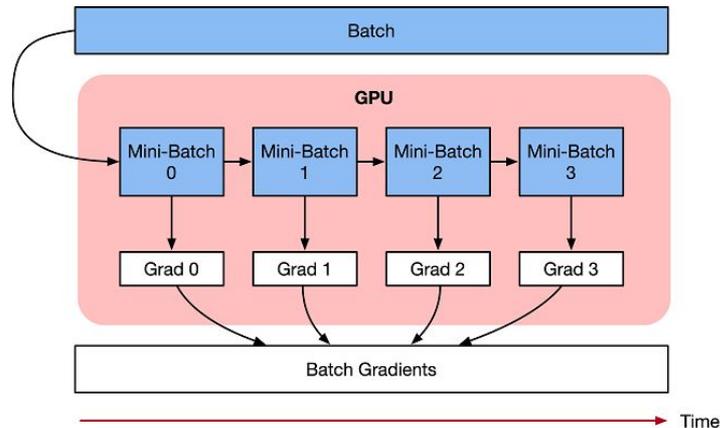  2. Repeat for **N** steps (accumulating gradients)
  3. optimizer.step() once.

**Trade-off:**
- ✅ **Memory:** Fits large logical batches on a single consumer GPU
- ❌ **Speed:** Sequential processing means it takes **N** times longer

Practical TIP: Let HF Accelerate handle accumulation for you!
→ Performing gradient accumulation with Accelerate

```
  from accelerate import Accelerator
- accelerator = Accelerator()
+ accelerator = Accelerator(gradient_accumulation_steps=2)
```



```
# Simulating batch_size = 64 with micro_batch = 8
accumulation_steps = 8

for i, batch in enumerate(dataloader):
    outputs = model(batch)
    loss = outputs.loss / accumulation_steps # Normalize!
    loss.backward() # Grads accumulate by default

    if (i + 1) % accumulation_steps == 0:
        optimizer.step() # Update once per 64 samples
        optimizer.zero_grad()
```

# Scaling Across GPUs for Speed:
# **Distributed Data Parallel (DDP)**

**Concept:** "Parallel Gradient Accumulation"

**Mechanism:**
- **Replication:** Every GPU holds an identical copy of the *entire* model
- **Data Splitting:** The batch is split across GPUs (e.g., GPU 0 gets frames 1-4, GPU 1 gets frames 5-8)
- **Sync:** Gradients are averaged across GPUs before the update step
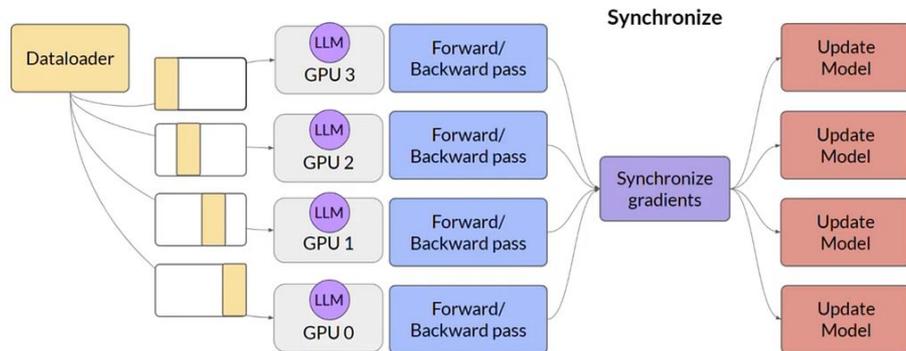
✅ **Pros:**
- **Speed:** Linear speedup (2 GPUs ≈ 2x faster)
- **Simplicity:** Minimal communication overhead

❌ **Cons:**
- **Memory Wall:** The *entire model* must fit on one GPU
- *Failure Mode:* You cannot train a 70B model on a 24GB card, no matter how many cards you have



Distributed Data Parallel (DDP)

# Scaling Across GPUs for *Size*:
# **Fully-Sharded Data Parallel (FSDP) // DeepSpeed ZeRO**

**Concept:** "Model Sharding"

**The Problem:** *What if the model itself is bigger than your VRAM?*
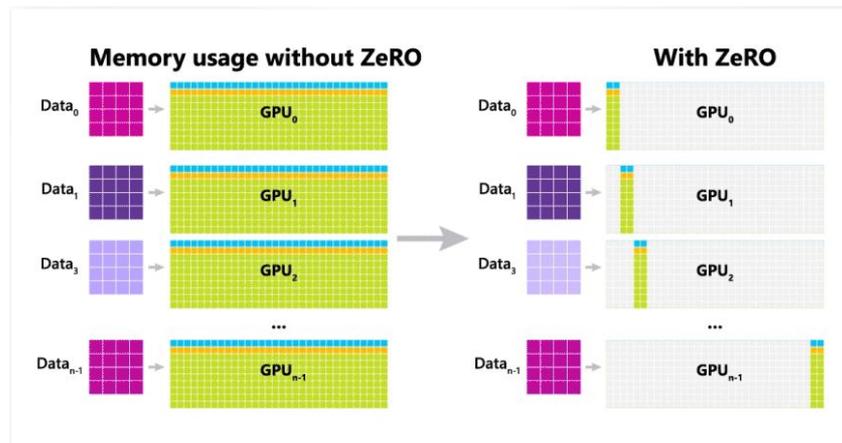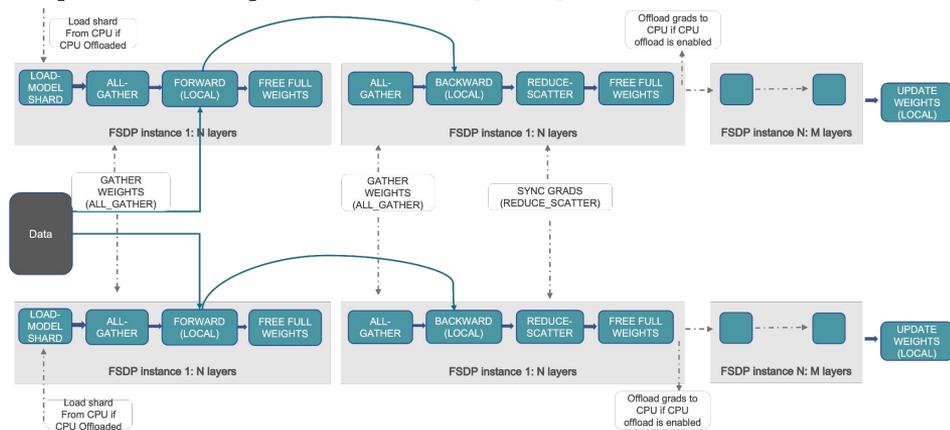
**Solution (ZeRO Stage 3):**

- **Shard Everything:** Split the Parameters, Gradients, and Optimizer States across all GPUs
- **On-Demand Fetching:** When GPU 0 needs Layer 1 to compute, it fetches those shards from other GPUs, computes, and discards them

✅ **Pros:**
- **Massive Scale:** Pools VRAM. 4x 24GB GPUs = ~96GB of usable Model Memory
- **Capability:** Enables fine-tuning 70B+ models on consumer clusters

❌ **Cons:**
- **Communication:** Heavy network traffic (constant swapping of shards)

# Shrinking the Model:
# **Quantization**

**Concept:**

- Represent weights with lower precision to save VRAM
- **FP16 (Half Precision):** 2 bytes/param (Standard Inference)
- **INT4 (Quantized):** 0.5 bytes/param (Extreme Efficiency)

**Example Impact (7B Model):**

- **FP16:** ~14GB VRAM (Requires A100/A6000).
- **INT4:** ~4GB VRAM (Fits on consumer RTX cards).

**The Stack:**

- **bitsandbytes:** The standard library for 4-bit/8-bit loading in PyTorch
- **NF4 (Normal Float 4):** A data type optimized for neural network weights (preserves accuracy better than standard INT4)

**In Practice:** can't train effectively in 4-bit (gradients vanish), so freeze the 4-bit weights and train a small adapter on top (QLoRA) →

# Training on a Budget:
# Parameter Efficient Finetuning ([PEFT](#))



**Partial Freezing**

- Freeze Vision Encoder (Always).
- Freeze LLM / *some-most* layers
- Train only the Projector

**LoRA (Low-Rank Adaptation)**

- Don't update the full weight matrix W
- Train low-rank decomposition matrices A and B where ΔW = BA
- → Reduces trainable params by 10,000x

**QLoRA (Quantized LoRA)**

- Load the base LLM in 4-bit (NF4) **[frozen]**
- Attach LoRA adapters in 16-bit **[trainable]**
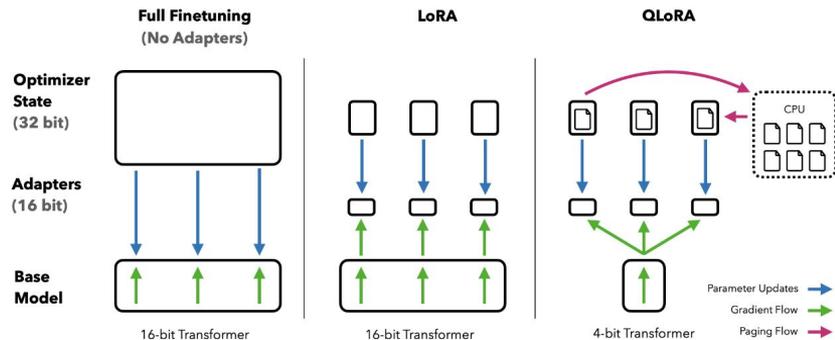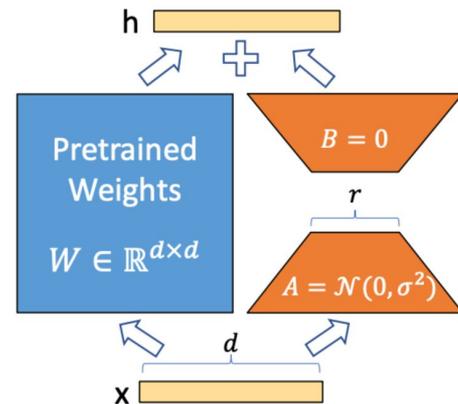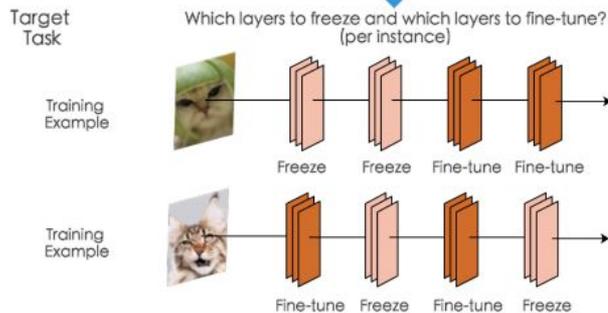- → Finetune a 70B model on a single 48GB GPU (A6000)



**Figure 1:** Different finetuning methods and their memory requirements. QLoRA improves over LoRA by quantizing the transformer model to 4-bit precision and using paged optimizers to handle memory spikes.

# V. Failure Modes

# *Where Video-LLMs still break* (non-exhaustive)

1. **Temporal confusion:** order, causality, "what changed?"

2. **Spatial grounding:** distances, collisions, occlusion, 3D ambiguity

3. **Long videos:** brute force attention scaling only goes so far

4. **Continual sensing:** tasks that require seeing *every* frame, fast responses

~all "Video" LLMs are secretly Image-LLMs in disguise!

→ frame subsampling is a persistent issue

# To be continued…

Week 7 (Mar 3) — **Tutorial 5**