# Task- YOLO Architectures,SOTA Analysis, Training Pipelines, and Videos Classification

(Week 5)

## YOLO & Darknet — notes, comparison, learning plan, and action checklist

Below is a compact, practical study guide can use to learn, summarize, diagram, and experiment with Darknet / YOLOv3 → YOLOv5 → YOLOv8 → YOLOv11. It includes (1) quick architecture summaries, (2) evolution highlights and comparison table, (3) where to read the primary sources and which files to inspect in Ultralytics repos, (4) suggested diagrams and how to draw them, and (5) an actionable experiment / training pipeline checklist. I looked up the canonical sources and the Ultralytics docs/repos to ensure accuracy. Ultralytics Docs+3ar5iv+3GitHub+3

## 1) Short architecture summaries

### Darknet (overview)

- **What it is:** Darknet is the original C-based framework used by the YOLO authors; includes several backbone networks used by YOLO variants (Darknet-19, Darknet-53).
- **Key point:** Darknet-53 (used by YOLOv3) is a 53-layer convolutional backbone designed for speed and competitive classification performance (compares to ResNet family but optimized for throughput). ar5iv

### YOLOv3 (Redmon & Farhadi, 2018)

- **Backbone:** Darknet-53 (residual-style conv blocks).
- **Neck / feature fusion:** Multi-scale predictions (predict on 3 feature-map scales — small/medium/large objects). Uses upsampling + route connections to fuse features.
- **Head / detection:** Each scale predicts bounding boxes via anchor boxes (k-means anchors), objectness score, and class probabilities. Uses logistic regression for objectness.

- **Loss & training notes:** Binary cross-entropy for objectness & class pred, sum-squared error for coordinates (older style). Multi-scale detection improved small-target recall. ar5iv

**Primary source:** *YOLOv3: An Incremental Improvement* (Redmon & Farhadi). arXiv

## YOLOv5 (Ultralytics — repo & docs)

- **Backbone:** Evolved designs in PyTorch — later versions used CSP-like ideas (CSP/Darknet variations appear across Ultralytics models). YOLOv5 is a practical PyTorch reimplementation that introduced many engineering/training conveniences (easy export, training scripts, augmentation pipelines). GitHub+1

- **Neck:** PANet-like feature fusion (path-aggregation) / SPP variants in different model sizes.

- **Head:** Anchor-based detection in early YOLOv5; later Ultralytics work moved toward anchor-free options in newer releases.

- **Why popular:** Excellent documentation, standardized training scripts, many pre-trained weights and export paths (ONNX, TFLite, TensorRT). GitHub

**Primary resources:** Ultralytics yolov5 repo & docs. GitHub+1

## YOLOv8 (Ultralytics)

- **Status:** Ultralytics' next-generation models and API improvements. No single formal academic paper — documented via Ultralytics docs and GitHub. YOLOv8 focuses on usability, modern PyTorch codebase, and incremental architecture/training improvements (CSP-inspired blocks, simplified heads, training/augmentation pipelines). Ultralytics Docs

- **Anchor-free options & efficiency:** Ultralytics moved toward anchor-free designs and more compact backbones in their newer model families; also emphasizes simplified training and export. Ultralytics Docs

## YOLOv11 (Ultralytics)

- **Status:** Latest Ultralytics release family (documentation page and GitHub discussion). As of Ultralytics docs there's a YOLOv11 model family; official paper may not exist — cite the repo/docs for specifics. Read the Ultralytics docs and release notes for architecture/training differences. Ultralytics Docs+1

## 2) Evolution highlights

| Area | YOLOv3 | YOLOv5 | YOLOv8 | YOLOv11 |
|---|---|---|---|---|
| Backbone | Darknet-53 (residual) | CSP-like / PyTorch backbones (varies) | CSPDarknet variants, lightweight options | further refined CSP / efficient blocks (see docs) |
| Neck | upsample + route (multi-scale) | PAN / SPP variants | PAN-like / optimized FPN | enhanced feature fusion / efficiency gains |
| Head | anchor-based multi-scale | anchor-based (early), modular export | anchor-free options available | anchor-free + training improvements |
| Training | basic augmentation, multi-scale | Rob. augmentations, TTA, export tools | easier pipelines, Ultralytics CLI | more SOTA training recipes (check release notes) |
| Papers | official paper | community/Ultralytics docs (no formal paper) | docs + repo | docs + repo/discussions |

*Takeaway:* YOLO evolved from a research prototype (Darknet + YOLOv3 paper) into a highly-engineered practical toolkit (YOLOv5+) where development emphasis shifted to efficient backbones (CSP), improved necks (PAN), anchor-free detection, and production-friendly tooling. ar5iv+1

## 3) Which files to inspect (Ultralytics repos) — exact hands-on steps

Clone the Ultralytics repos and inspect model definition / backbone modules:

```
# YOLOv5
git clone https://github.com/ultralytics/yolov5.git
cd yolov5
# inspect model definitions
# architectures are typically in models/ or models/common.py
ls models
sed -n '1,240p' models/common.py

# Ultralytics (new unified repo for v8+/v11)
git clone https://github.com/ultralytics/ultralytics.git
cd ultralytics
# models often live in ultralytics/models or ultralytics/yolo
ls ultralytics/models
```

Files and classes to look for:

- models/common.py, models/yolo.py (model blocks, CSP modules).

- models/backbone.py or cspdarknet.py (backbone implementations).

- models/heads.py or head-related functions (detection head, anchor/head logic).

- train.py, val.py — show training recipes and augmentation pipelines. GitHub+1

## 4) Diagrams — what to draw & how (Draw.io / diagrams.net steps)

### Diagrams to produce

1. **Darknet-53 block diagram** — show residual blocks and overall 53-layer flow (stem → conv blocks → classifier).

2. **YOLOv3 architecture** — backbone (Darknet-53), neck (feature upsample + route), 3-scale detection head.

3. **YOLOv5 / YOLOv8 architecture** — backbone (CSPDarknet / CSP modules), neck (PAN), head (detect layers). Annotate where SPP, PANet, CSP split/concat happen.

4. **Comparison diagram** — small side-by-side of backbone/necks/heads across the versions, highlighting CSP, anchor-free, and SPP changes.

### How to draw quickly

- Use Draw.io / diagrams.net. Start with 3 vertical columns (Backbone / Neck / Head) and draw colored boxes for each module (conv block, CSP block, PAN). Add arrows for data flow and annotate feature-map sizes (e.g., 80×80 → 40×40 → 20×20 for 320px input).

- Export PNG + SVG for repo. Include a short caption under each diagram.

## 5) Practical checklist for summarizing & comparing (deliverable)

1. **Read & cite**

   - YOLOv3 paper. arXiv

   - Ultralytics YOLOv5 docs & repo. GitHub

   - Ultralytics YOLOv8 docs. Ultralytics Docs

   - Ultralytics YOLOv11 docs / release notes & community discussion. Ultralytics Docs+1

2. **Inspect code**

- Clone repos and find models/ files (CSP modules, Detect head). [GitHub+1](#)

3. **Take notes** — for each version record:

   - Backbone blocks (e.g., residual conv, CSP split), parameter counts, typical receptive fields.

   - Neck (FPN / PAN / SPP), where features are fused.

   - Head (anchors vs anchor-free, how bounding boxes are parameterized).

   - Training tricks: augmentations, schedulers, loss terms (objectness vs classification), post-processing (NMS variant).

4. **Make diagrams** (one per version + one comparative).

5. **Write a short evolution summary** (200–400 words) with a table listing differences (backbone, neck, head, anchor vs anchor-free, SOTA notes).

# 6) Recommended focused experiments / probing

- **Inspect receptive fields & FLOPs** of each backbone with a profiler (torchinfo / thop). Measure inference FPS on a GPU.

- **Ablation**: run small experiments toggling anchor vs anchor-free head (or SPP on/off) to see mAP vs speed tradeoffs.

- **Visualize feature maps** from the backbone (intermediate layers) and the outputs of PAN to understand what scales detect which objects.

- **Reproduce a small COCO subset** training with YOLOv5 and YOLOv8 to compare quality / speed. Use Ultralytics train.py (they provide scripts) and export best models. [GitHub+1](#)

# 7) Quick-reading path & prioritized resources

1. **YOLOv3 paper (read thoroughly)** — understand Darknet-53, multi-scale detection, and anchors. [arXiv](#)

2. **Ultralytics YOLOv5 GitHub + docs** — practical implementation and engineering choices; read models/common.py, models/experimental.py, and training docs. [GitHub+1](#)

3. **Ultralytics YOLOv8 docs** — read the model overview and migration notes (anchor-free info). [Ultralytics Docs](#)

4. **Ultralytics YOLOv11 docs / release notes** — read official docs and community discussion for latest changes. [Ultralytics Docs+1](#)

## 8) Short annotated bibliography (links)

- YOLOv3 paper — Redmon & Farhadi (2018). arXiv
- Ultralytics YOLOv5 repo & docs. GitHub+1
- Ultralytics YOLOv8 docs. Ultralytics Docs
- Ultralytics YOLOv11 docs & community thread. Ultralytics Docs+1
- PapersWithCode / benchmarks for YOLOv3 variants. paperswithcode.com

## 9) Deliverables

- A concise **Markdown summary** per version (YOLOv3 / YOLOv5 / YOLOv8 / YOLOv11) with pros/cons and code-file pointers.
- A **comparison table** (CSV/Markdown) listing backbone, neck, head, anchors, typical FLOPs / params (estimates).
- A **Draw.io diagram** (XML) template for one architecture (e.g., YOLOv3) you can open in diagrams.net.
- A small **notebook** that clones Ultralytics repo, parses model config files, and prints a readable summary of blocks and channels.
- A **bash snippet** to clone repos and list model files as shown above.

# Executive summary

- Modern SOTA for **Image Classification** is dominated by large image–text / multimodal backbones (e.g., CoCa) or large ViT / ConvNeXt variants; CoCa (finetuned) holds a reported **~91.0% top-1** on ImageNet (finetuned). Papers with Code+1
- For **Object Detection / Instance Segmentation** the top entries on COCO (as of mid-2025) are transformer-based query / DETR-style models and large foundation backbones (e.g., **Co-DETR** reporting **~66.0 box AP** on COCO test-dev; InternImage variants are very close). CVF Open Access+1
- **YOLO family (Ultralytics)** continues to focus on speed/efficiency and real-world deployment (YOLO11 variants trade some absolute AP for much smaller parameters / faster inference). Example: **YOLO11m ≈ 51.5 mAP (COCO)** with ~20M params — far smaller and faster than the large DETR/ViT detectors but lower absolute AP. Ultralytics Docs
- For **segmentation**, unified query/mask systems such as **Mask2Former / MaskDINO** remain highly competitive on mask-AP metrics; foundation models for

promptable segmentation (SAM → **SAM 2**) provide powerful zero-shot / interactive segmentation (esp. images & video). arXiv+1

- **Tradeoff**: SOTA detectors → higher absolute AP but larger models, longer train / inference cost; YOLO → smaller, lower latency, very practical for production and edge. Use the model family that matches your constraints (accuracy vs. latency vs. memory).

# 1) Darknet & YOLO architectures

## Darknet (legacy backbone)

- **Darknet-53** (original YOLOv3 backbone) is a 53-layer convolutional feature extractor designed for real-time detection; it uses residual connections and successive downsampling to produce multi-scale features used by the YOLO head. It was the backbone in YOLOv3. arXiv+1

## YOLOv3 (Redmon & Farhadi)

- Introduced multi-scale prediction (three feature maps) and used Darknet-53; emphasis on speed with competitive mAP in 2018. Architecture = **backbone (Darknet-53) → neck (feature pyramids) → head (prediction at scales)**. arXiv

## YOLOv4 / YOLOv5 / Late Ultralytics variants — practical changes

- From YOLOv4 onward the community and Ultralytics introduced many practical improvements: CSP modules, PAN/FPN necks, improved augmentation, and training recipes. YOLOv5 (and later Ultralytics releases) are strong on usability and inference pipelines (PyTorch toolchain + export). GitHub+1

## YOLOv8 → YOLOv11 (Ultralytics): modern Ultralytics innovations

- **Anchor-free split head** (decoupled classification, box, mask outputs), CSP / CSPDarknet style improvements in backbone, new neck modules (e.g., C2f), and a strong emphasis on deployment/export (ONNX, TensorRT). YOLO11 emphasizes efficiency & multi-task variants (detection, seg, pose) with a family (n/s/m/l/x). Ultralytics documentation provides per-variant mAP, params, and speed. Ultralytics Docs+1

## Key architectural components (across YOLO family)

- **Backbone** — feature extractor (Darknet → CSPDarknet → custom Ultralytics backbones).

- **Neck** — multi-scale feature aggregator (FPN, PAN, C2f, etc.).

- **Head** — final prediction module (anchor-based historically; anchor-free split head in modern Ultralytics).

- **Training/recipes** — extensive augmentations (mosaic, mixup/cutmix variations), multi-scale training, LR schedules, and pruning/quantization for deployment.

# 2) SOTA model analysis (classification, detection, segmentation) — short list, metrics & reasoning

I prioritized high-impact leaderboard numbers (ImageNet for classification; COCO for detection/segmentation) and recent foundation models/papers. Below are the top contenders and the numbers reported in public sources.

## Image classification (ImageNet top-1)

- **CoCa (Contrastive Captioners)** — reported **91.0% top-1 (finetuned)** on ImageNet (paper + PapersWithCode entry). This model is a multimodal image-text foundation model and is one of the highest reported single-model ImageNet top-1s (large parameter counts reported). arXiv+1

- Other top performers (close secondaries) include very large ViT/ConvNeXt variants, "model soups" and large multimodal/vision foundation models; these often trade compute & parameter count for the last few tenths of percent. Papers with Code

## Object detection (COCO box AP — test-dev / val)

- **Co-DETR (DETR with collaborative hybrid assignment)** — **≈66.0 box AP (COCO test-dev)** for ViT-Large variant (reported in paper/readme). This is among the top node on the COCO detection leaderboard in recent listings. CVF Open Access+1

- **InternImage-H** — reported **≈65.4 AP** (InternImage family achieved strong detection records using deformable convolutions / DCNv3 variants). CVF Open Access

- **DINO / MaskDINO** family — strong DETR-family competitors; DINO reported **~63.2/63.3 AP** (val/test-dev) in earlier SOTA cycles; MaskDINO extended DINO into seg and improved unified mask APs. Hugging Face+1

## Instance / panoptic segmentation (COCO mask AP / PQ)

- **Mask2Former / MaskDINO / Mask Frozen-DETR / Co-DETR** variants report top mask APs (Mask2Former provided a big leap when introduced; MaskDINO reports ~54.7 mask AP on some leaderboards). Foundation models (SAM → SAM 2) changed the application landscape: excellent promptable segmentation and video support, but task-specific leaderboard numbers depend on fine-tuning / evaluation protocol. arXiv+2CVF Open Access+2

# 3) YOLO vs SOTA detectors — concrete comparison (accuracy vs efficiency)

**Example numeric comparisons (representative figures from vendor/paper pages):**

- **YOLO11m**: COCO box mAP ≈ **51.5** (params ≈ **20.1M** — small, fast family). Ultralytics Docs

- **YOLO11x**: COCO box mAP ≈ **54.7** (params larger ≈ 56.9M). Ultralytics Docs

- **Co-DETR (ViT-L)**: COCO box AP ≈ **66.0** (params ≈ **304M** for ViT-Large DETR variant). CVF Open Access

**Interpretation / tradeoffs**

- Big DETR/ViT / foundation backbones reach substantially higher COCO AP (often +10–15 AP over YOLO mid-size models), but at much higher parameter counts, training compute, and inference cost. CVF Open Access+1

- YOLO family (Ultralytics) is optimized for practical deployment: faster inference, smaller model sizes, multi-task variants (segmentation/pose), and export pipelines (ONNX/TensorRT). If you require real-time constraints (edge/embedded) YOLO usually wins; if you need absolute top AP and compute is available, DETR/large backbones are preferred. Ultralytics Docs+1

# 4) Training pipeline practical notes

## Common items for high performing detectors

1. **Large pretraining + fine-tuning** — transfer learning from strong backbones (ImageNet / multi-modal pretraining) is standard for top AP. arXiv+1

2. **Data augmentation** — mosaic, mixup/cutmix, random scale & flip, color jitter, AutoAugment/RandAugment where useful (Ultralytics & community recipes). Ultralytics Docs

3. **LR scheduling** — cosine annealing (with restarts if desired) or one-cycle; longer schedules and staged LR drops help DETR-style models.

4. **Mixed precision (AMP)** and distributed training to speed up large runs.

5. **Loss / matching** — DETR style uses Hungarian matching (one-to-one), improved denoising and other tricks (DINO, MaskDINO) to stabilize and accelerate training. CVF Open Access+1

## For YOLO / real-time families

- Multi-scale training, knowledge distillation (for smaller models), quantization / pruning and TensorRT export are common optimizations. Ultralytics provides recipes & per-variant benchmarks. Ultralytics Docs

## For segmentation (mask quality)

- High-resolution feature maps, mask heads (pixel-embedding + dot-product), mask-specific decoders and mask loss weighting; query-based mask predictors (MaskDINO / Mask2Former) are strong. arXiv+1

# 5) Video classification & video pipelines

- **Backbones / SOTA models**: VideoMAE (masked autoencoder for video), Video Swin / MViT / TimeSformer, and large variants (VideoMAE V2-g reported very strong Kinetics results). On Kinetics/SSv2 leaderboards, these transformer-style or masked-pretraining models occupy top slots. Example: **VideoMAE V2-g ≈ 88.5** on Kinetics-400 in some listings. Papers with Code+1

- **Pipelines**: temporal sampling (sparse or dense), tube masking for pretraining, optionally optical-flow streams for motion, and frame/clip augmentation. Use pretraining on large video corpora when possible; finetune on Kinetics / SSv2. GitHub+1

# 6) Real-time / near-real-time DETRs: DETRs Beat YOLOs on real-time?

- Work such as **RT-DETR** (Real-Time DETR) and its followups show that transformer-based detectors can be engineered for real-time performance and can in some cases match/beat YOLOs on the speed/accuracy frontier by careful hybrid encoder design & engineering (TensorRT etc.). But the hardware + implementation choices matter a lot (accelerated backends, TensorRT/ONNX). arXiv+1

# 7) Recommendations

## If goal is **maximum COCO AP**:

- Use a DETR or query-based detector with a strong backbone (ViT-Large / InternImage-H / pretrained DINO family), long training, large batch / multi-GPU, optional Objects365 or web data pretraining, MaskDINO/Co-DETR choices. Expect high compute & memory. CVF Open Access+1

## If goal is **production / real-time**:

- Use current Ultralytics YOLO11 family (pick model size per latency budget). Use quantization / TensorRT export; apply Ultralytics best practices and augmentation recipes. Ultralytics docs provide per-variant mAP / params / FLOPs to guide selection. Ultralytics Docs

## If goal is **segmentation or interactive segmentation**:

- Use Mask2Former / MaskDINO / SAM 2 depending on task: training for mask AP vs promptable zero-shot segmentation. For video promptable work, SAM 2 offers improved video support. arXiv+1

## Video classification experiments:

- Use VideoMAE or Video Swin backbones; pretrain (masked autoencoder) and then finetune on Kinetics / SSv2; evaluate top-1/top-5 & perplexity where appropriate. GitHub+1

# 8) Suggested deliverables / repo layout to support reproducible evaluation

Structure the GitHub repo (pytorch-week3) like this (short listing):

```
pytorch-week3/
 README.md
 code/
  cls/        # ResNet-18 CIFAR pipeline (from-scratch code)
  mt/         # Transformer toy translation code & notebook
  det/        # Detection experiments: configs for YOLO11 / RT-DETR / Co-DETR (where allowed)
  seg/        # Segmentation experiments
  video/      # Video classification pipelines
 runs/
  cls/ (figures: curves_cls.png, confusion_matrix.png, preds_grid.png, miscls_grid.png, gradcam_*.png)
```

```
 mt/  (curves_mt.png, attention_layerL_headH.png, masks_demo.png, decodes_table.png, bleu_report.png)
 det/
 seg/
 video/
notebooks/
 cls_visualize.ipynb   # plots & Grad-CAM displays
 mt_interactive.ipynb  # training / attention heatmaps (sacrebleu)
scripts/
 run.sh  # baseline, mixup, cutmix, long-run commands
report/
 report.md
 onepage_report.pdf
docs/
 architecture_diagrams.drawio
```

# 9) Concrete citations

- YOLOv3 — Redmon & Farhadi, *YOLOv3: An Incremental Improvement* (arXiv / paper). arXiv

- Ultralytics YOLO11 docs (per-variant mAP, speed & params). Ultralytics Docs

- CoCa (Contrastive Captioners) — paper + ImageNet entry (91.0% top-1 reported). arXiv+1

- Co-DETR — ICCV/ArXiv / code and COCO results (≈66.0 AP, ViT-L). CVF Open Access+1

- InternImage (InternImage-H ≈65.4 mAP on COCO) — InternImage paper & repo. CVF Open Access

- Mask2Former (masked-attention Mask Transformer for universal segmentation). arXiv

- SAM 2 (Segment Anything Model 2) — Meta (paper / tech blog). arXiv+1

- DINO / MaskDINO — DETR improvements and mask metrics. CVF Open Access+1

- VideoMAE / Video classification SOTA listing (Kinetics / Something-Something metrics). GitHub+1

# 10) Next steps & experiments

1. **Small reproducible DETR vs YOLO benchmark** (single GPU): train YOLO11s on COCO8 (toy subset) and a small DETR variant (RT-DETR small) — measure

throughput (FPS) and AP after 50–100 epochs; collect export timings (ONNX / TRT). Use Ultralytics guide & RT-DETR repo. Ultralytics Docs+1

2. **ResNet / CIFAR run**: implement ResNet-18 from primitives, train baseline → add mixup / cutmix → warmup LR → compare curves & confusion matrix. Save best .pt and log to TensorBoard. (This matches your earlier spec.)

3. **Transformer toy MT**: implement minimal Transformer encoder-decoder (from primitives), build an expanded toy corpus (questions, negations, varied word orders), train to stable BLEU ≥15; produce attention heatmaps & sacrebleu report.

4. **Video classification**: fine-tune a VideoMAE / Video Swin small variant on Kinetics-400 (or a tiny subset for experimentation); produce top-1 curves, confusion matrices and some sample predictions. GitHub+1

# Appendix — short notes on reproducibility & benchmarking

- When comparing models, always record: dataset split (val/test), input image size, single- vs multi-scale evaluation, TTA, and whether pretraining used extra private data (many SOTA results rely on large pretraining corpora). Leaderboard numbers can be sensitive to these details — always match evaluation protocol. CVF Open Access+1

**A drop-in** report.md **template** with clear section headings, placeholders for images/figures, tables for metrics, and spots for your own insights:

# Deep Learning Experiments Report

## 1. Introduction
- **Objective:**
  _(state the experiment's goal, e.g., CIFAR-10 classification, COCO detection, toy MT, etc.)_
- **Models Compared:**
  _(list models: ResNet-18, YOLO11, Co-DETR, Transformer toy, etc.)_
- **Datasets Used:**
  _(CIFAR-10, COCO, WMT toy dataset, Kinetics subset, etc.)_

---

## 2. Methods

### 2.1 Architectures
- **Backbones:** _(e.g., ResNet-18, Darknet-53, ViT-L, etc.)_

- **Heads:** _(classification/detection/segmentation heads used)_
- **Variants:** _(baseline, MixUp, CutMix, label smoothing, cosine LR, etc.)_

### 2.2 Training Setup
- **Hardware:** _(GPU/CPU details)_
- **Hyperparameters:** _(batch size, LR, optimizer, scheduler, epochs)_
- **Augmentations:** _(MixUp, CutMix, RandAugment, Cutout)_

---

## 3. Results

### 3.1 Classification (CIFAR-10 / ImageNet subset)
- **Learning Curves:**
  ![Train/Val Curves](runs/cls/curves_cls.png)
  *Figure 1. Training and validation accuracy/loss curves.*

- **Confusion Matrix:**
  ![Confusion Matrix](runs/cls/confusion_matrix.png)
  *Figure 2. Confusion matrix on validation set.*

- **Sample Predictions (Correct vs Misclassified):**
  ![Correct vs Misclassified](runs/cls/preds_grid.png)
  *Figure 3. Examples of predictions with true/false labels.*

- **Grad-CAM Overlays:**
  ![Grad-CAM](runs/cls/gradcam_sample.png)
  *Figure 4. Grad-CAM visualization highlighting important regions.*

- **Table of Results:**

| Model     | Augmentation | Top-1 Acc (%) | Params (M) | Notes |
|-----------|--------------|---------------|------------|-------|
| ResNet-18 | Baseline     | xx.x          | xx.x       | –     |
| ResNet-18 | MixUp        | xx.x          | xx.x       | better stability |
| ResNet-18 | CutMix       | xx.x          | xx.x       | improved generalization |

---

### 3.2 Object Detection (COCO)

- **PR Curves:**
  ![PR Curves](runs/det/pr_curves.png)
  *Figure 5. Precision-Recall curves.*

- **mAP Table:**

| Model     | Backbone    | mAP@0.5:0.95 | Params (M) | FPS | Notes |

```
|-------------|------------|-------------|-----------|-----|-------|
| YOLO11m    | CSPDarknet | xx.x        | xx.x      | xx  | fast |
| YOLO11x    | CSPDarknet | xx.x        | xx.x      | xx  | more accurate |
| Co-DETR (L) | ViT-Large  | xx.x        | xx.x      | xx  | highest AP |
```

- **Sample Detections:**
  ![Detections](runs/det/sample_dets.png)
  *Figure 6. Example detections with bounding boxes.*

---

### 3.3 Segmentation

- **Mask Quality:**
  ![Segmentation Masks](runs/seg/mask_samples.png)
  *Figure 7. Predicted vs ground-truth masks.*

- **Metrics Table:**

| Model       | Backbone | Mask AP | PQ | Notes |
|-------------|----------|---------|----|-------|
| YOLO11-seg  | CSPDarknet | xx.x  | xx | fast |
| Mask2Former | Swin-L   | xx.x    | xx | best overall |
| SAM 2       | ViT-Huge | n/a (promptable) | – | zero-shot segmentation |

---

### 3.4 Machine Translation (Toy Transformer)

- **Learning Curves:**
  ![MT Curves](runs/mt/curves_mt.png)

- **Attention Heatmaps:**
  ![Attention](runs/mt/attn_layer3_head2.png)

- **Sample Translations:**

| Source | Prediction | Reference | Notes |
|--------|------------|-----------|-------|
| "She is not coming." | "Elle ne vient pas." | "Elle n'arrive pas." | acceptable |

- **BLEU Scores:**

| Epoch | BLEU |
|-------|------|
| 10    | xx.x |
| 20    | xx.x |

---

### 3.5 Video Classification

- **Confusion Matrix:**
  ![Video Confusion](runs/video/confusion_matrix.png)

- **Top-K Accuracy Table:**

  | Model       | Dataset     | Top-1 (%) | Top-5 (%) | Notes |
  |-------------|-------------|-----------|-----------|-------|
  | VideoMAE-S  | Kinetics-400 | xx.x     | xx.x      | baseline |
  | SwinV2-B    | Kinetics-400 | xx.x     | xx.x      | better accuracy |

---

## 4. Discussion
- **Key Insights:**
  _(compare models, note overfitting/generalization, augmentation impacts, efficiency vs accuracy tradeoffs, etc.)_

- **YOLO vs SOTA:**
  _(summarize differences in AP vs speed, production-readiness vs research SOTA)_

- **Limitations:**
  _(what didn't work well, e.g., unstable BLEU, segmentation runtime heavy)_

---

## 5. Conclusion
- **Summary of Results**
- **Best Performing Models**
- **Future Work / Improvements**

---

## References
_(insert the reference list you prepared earlier)_

# Training Pipeline & Data Preparation – Plan

## 1. Repository Structure

yolo-training-pipeline/
|
├── code/

```
|   ├── data_prep.py      # dataset download/format conversion
|   ├── train_yolo.py     # training script (PyTorch/Ultralytics API)
|   ├── utils.py          # helper functions (metrics, plotting)
|   └── __init__.py
|
├── configs/
|   ├── coco.yaml         # dataset config (COCO format)
|   ├── voc.yaml          # dataset config (Pascal VOC format)
|   └── hyp.yaml          # training hyperparameters
|
├── runs/
|   └── exp1/             # logs, tensorboard, model weights, plots
|
├── notebooks/
|   └── visualize_data.ipynb  # check augmentations, bounding boxes
|
├── report/
|   └── report.md         # results, figures, discussion
|
├── run.sh                # bash commands for training runs
├── README.md
└── requirements.txt
```

## 2. Data Preparation (data_prep.py)

- Download datasets: COCO or VOC.

- Convert to YOLO format (txt per image: class x_center y_center w h).

- Apply **Albumentations** or **torchvision.transforms** for augmentation.

- Mosaic augmentation (4-image merge) for YOLO-style training.

- Document preprocessing in the README.

## 3. Training Script (train_yolo.py)

- Use Ultralytics YOLOv8 API (from ultralytics import YOLO) or custom PyTorch training loop.

- Support **optimizers**: SGD (with momentum, warmup) / AdamW.

- Loss: BCE + CIoU (standard YOLO loss).

- CLI args for:

    - epochs, batch size, learning rate, optimizer

- augmentation flags (mosaic, mixup, hsv shift)
- resume from checkpoint
- Save **best.pt** (best model checkpoint).
- Export results: confusion matrix, PR curves, mAP table.

## 4. Training Run Example (run.sh)

```bash
#!/bin/bash

# Baseline YOLOv8n, VOC dataset
yolo detect train data=configs/voc.yaml model=yolov8n.pt epochs=20 imgsz=640 batch=16 \
  optimizer=SGD lr0=0.01 momentum=0.937 weight_decay=0.0005 \
  name=exp_voc_baseline

# With Mosaic + MixUp
yolo detect train data=configs/voc.yaml model=yolov8n.pt epochs=20 imgsz=640 batch=16 \
  augment=True mosaic=1.0 mixup=0.2 \
  name=exp_voc_augmented

# COCO subset run
yolo detect train data=configs/coco.yaml model=yolov8n.pt epochs=10 imgsz=640 batch=32 \
  name=exp_coco_subset
```

## 5. Visualizations & Logging

- TensorBoard / wandb logging (loss curves, PR curves, confusion matrix).
- Notebook visualize_data.ipynb to preview augmentations & bounding boxes.
- Figures saved into runs/expX/.

## 6. Deliverables

- **Training scripts** (train_yolo.py, data_prep.py).
- **Configs** for VOC/COCO.
- **run.sh** with tuned runs.
- **Report.md** with:
  - Dataset details
  - Hyperparameters used

- Results: mAP, precision, recall

- Figures: curves, confusion matrices, sample detections

- Insights on augmentation & optimizer choices.

**A ready-to-drop** train_yolo.py script that uses the Ultralytics API, adds augmentation toggles, checkpoint saving, TensorBoard logging, and CLI arguments for flexibility.

```python
#!/usr/bin/env python3
"""
train_yolo.py
-------------
Train YOLOv8 models on VOC or COCO datasets using Ultralytics API.
Includes:
- CLI arguments for flexibility
- Augmentation flags (mosaic, mixup, hsv)
- Optimizer options (SGD, AdamW)
- Logging (TensorBoard, CSV, JSON)
- Auto checkpoint saving (best.pt, last.pt)
"""

import argparse
from ultralytics import YOLO

def parse_args():
    parser = argparse.ArgumentParser(description="YOLOv8 Training Script")

    # Model & Data
    parser.add_argument("--model", type=str, default="yolov8n.pt",
            help="Pretrained YOLOv8 model to fine-tune (e.g., yolov8n.pt, yolov8s.pt)")
    parser.add_argument("--data", type=str, required=True,
            help="Path to dataset config YAML (VOC/COCO in YOLO format)")
    parser.add_argument("--imgsz", type=int, default=640, help="Image size (default: 640)")

    # Training
    parser.add_argument("--epochs", type=int, default=20, help="Number of epochs")
    parser.add_argument("--batch", type=int, default=16, help="Batch size")
    parser.add_argument("--optimizer", type=str, default="SGD", choices=["SGD", "AdamW"],
            help="Optimizer to use")
    parser.add_argument("--lr0", type=float, default=0.01, help="Initial learning rate")
    parser.add_argument("--momentum", type=float, default=0.937, help="Momentum for SGD")
    parser.add_argument("--weight_decay", type=float, default=0.0005, help="Weight decay")

    # Augmentations
    parser.add_argument("--augment", action="store_true", help="Enable default YOLO augmentations")
    parser.add_argument("--mosaic", type=float, default=1.0, help="Mosaic augmentation probability")
    parser.add_argument("--mixup", type=float, default=0.0, help="MixUp augmentation probability")
```

```python
    parser.add_argument("--hsv", type=float, default=0.015, help="HSV augmentation gain")

    # Checkpointing & Logging
    parser.add_argument("--project", type=str, default="runs/train", help="Project dir")
    parser.add_argument("--name", type=str, default="exp", help="Run name")
    parser.add_argument("--exist_ok", action="store_true", help="Allow existing project/name")
    parser.add_argument("--resume", action="store_true", help="Resume training from last checkpoint")
    parser.add_argument("--device", type=str, default="",
                help="CUDA device (e.g., 0,1) or 'cpu'")

    return parser.parse_args()


def main():
    args = parse_args()

    # Load model
    model = YOLO(args.model)

    # Training configuration
    results = model.train(
        data=args.data,
        imgsz=args.imgsz,
        epochs=args.epochs,
        batch=args.batch,
        optimizer=args.optimizer,
        lr0=args.lr0,
        momentum=args.momentum,
        weight_decay=args.weight_decay,
        augment=args.augment,
        mosaic=args.mosaic,
        mixup=args.mixup,
        hsv_h=args.hsv,
        hsv_s=args.hsv,
        hsv_v=args.hsv,
        project=args.project,
        name=args.name,
        exist_ok=args.exist_ok,
        resume=args.resume,
        device=args.device,
        save=True,          # save checkpoints
        save_period=5,      # save every 5 epochs
        cache=False,        # set True if enough RAM
        verbose=True,
    )

    print("  Training complete!")
    print(f"  Best model saved to: {results.save_dir}/weights/best.pt")
```

```
print(f" Last checkpoint: {results.save_dir}/weights/last.pt")


if __name__ == "__main__":
  main()
```

## ---- How to Run --------

Example baseline run:

```
python code/train_yolo.py \
 --model yolov8n.pt \
 --data configs/voc.yaml \
 --epochs 20 --batch 16 --optimizer SGD \
 --project runs/train --name exp_voc_baseline
```

With augmentations:

```
python code/train_yolo.py \
 --model yolov8n.pt \
 --data configs/voc.yaml \
 --epochs 20 --batch 16 --optimizer SGD \
 --augment --mosaic 1.0 --mixup 0.2 \
 --project runs/train --name exp_voc_aug
```

Resume from checkpoint:

```
python code/train_yolo.py --resume --project runs/train --name exp_voc_baseline
```

## ------ Research Notes: Video Classification Models -------

### Classic Approaches

- **CNN + RNN (e.g., ResNet + LSTM)**
    - CNN extracts per-frame spatial features.
    - RNN (LSTM/GRU) aggregates temporal info.
    - Pros: Simple, interpretable, lightweight.
    - Cons: Temporal modeling is limited, training can be slow.

### 3D CNNs

- **C3D / I3D (Inflated 3D ConvNet)**
    - Extend 2D conv filters to 3D (spatio-temporal).

- I3D pre-trained on Kinetics is strong baseline.

- Pros: Captures spatio-temporal info directly.

- Cons: Heavy computation, high memory usage.

### Transformer-based

- **VideoMAE, TimeSformer, ViViT**

  - Treat video frames as patches (like ViT), apply temporal attention.

  - Pros: SOTA performance, flexible, scalable.

  - Cons: Requires lots of data & compute.

## Baseline to Implement

Since spec says **"implement at least one baseline model"**, a good starting point is:

👉 **CNN + LSTM Baseline**
(using ResNet backbone + LSTM temporal model)

This is light, easy to train on **UCF-101 mini** or even **Kinetics-400 (subset)**.

## Suggested Repo Structure Update

```
pytorch-week3/
|
├── code/
|   ├── train_resnet.py
|   ├── train_transformer.py
|   ├── train_yolo.py
|   ├── train_video.py       <-- NEW (CNN+LSTM baseline)
|   └── models/
|      ├── resnet.py
|      ├── transformer.py
|      └── video_cnn_lstm.py <-- NEW
|
├── runs/
|   ├── cls/
|   ├── mt/
|   └── video/          <-- NEW (logs, figs, confusion matrix)
|
├── report/
|   ├── report.md
|   └── onepage_visual.md
```

# Draft video_cnn_lstm.py (Baseline Model)

```python
import torch
import torch.nn as nn
import torchvision.models as models


class CNNLSTMVideoClassifier(nn.Module):
    def __init__(self, num_classes: int = 101, hidden_dim: int = 256, num_layers: int = 1):
        super().__init__()

        # Pretrained ResNet for frame-level feature extraction
        resnet = models.resnet18(pretrained=True)
        modules = list(resnet.children())[:-1]  # remove fc
        self.cnn = nn.Sequential(*modules)
        self.feature_dim = resnet.fc.in_features

        # Freeze CNN backbone initially (can unfreeze later)
        for p in self.cnn.parameters():
            p.requires_grad = False

        # LSTM for temporal modeling
        self.lstm = nn.LSTM(
            input_size=self.feature_dim,
            hidden_size=hidden_dim,
            num_layers=num_layers,
            batch_first=True
        )

        # Final classifier
        self.fc = nn.Linear(hidden_dim, num_classes)

    def forward(self, x):
        """
        x: [B, T, C, H, W]  (batch, time, channels, height, width)
        """
        B, T, C, H, W = x.shape
        x = x.view(B * T, C, H, W)

        # Frame-level features
        with torch.no_grad():
            feats = self.cnn(x).view(B, T, -1)

        # LSTM
        out, _ = self.lstm(feats)  # [B, T, hidden_dim]
```

```
    # Last timestep
    out = out[:, -1, :]

    return self.fc(out)
```

## Draft train_video.py

This script:

- Loads UCF-101 frames.
- Uses CNN+LSTM model.
- Logs metrics & saves confusion matrix.

```python
#!/usr/bin/env python3
import argparse
import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
from torchvision import transforms
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
import matplotlib.pyplot as plt
import os

from models.video_cnn_lstm import CNNLSTMVideoClassifier
from datasets.ucf101 import UCF101Dataset  # <-- you'll need a simple dataset wrapper


def parse_args():
    parser = argparse.ArgumentParser(description="Train CNN+LSTM Video Classifier")
    parser.add_argument("--data_root", type=str, required=True, help="Path to UCF-101 frames")
    parser.add_argument("--epochs", type=int, default=10)
    parser.add_argument("--batch", type=int, default=8)
    parser.add_argument("--lr", type=float, default=1e-3)
    parser.add_argument("--num_classes", type=int, default=101)
    parser.add_argument("--device", type=str, default="cuda" if torch.cuda.is_available() else "cpu")
    parser.add_argument("--out_dir", type=str, default="runs/video")
    return parser.parse_args()


def main():
    args = parse_args()
    os.makedirs(args.out_dir, exist_ok=True)

    # Dataset + Dataloader
```

```
transform = transforms.Compose([
    transforms.Resize((112, 112)),
    transforms.ToTensor()
])
train_set = UCF101Dataset(args.data_root, split="train", transform=transform, frames_per_clip=16)
val_set = UCF101Dataset(args.data_root, split="val", transform=transform, frames_per_clip=16)

train_loader = DataLoader(train_set, batch_size=args.batch, shuffle=True)
val_loader = DataLoader(val_set, batch_size=args.batch, shuffle=False)

# Model
model = CNNLSTMVideoClassifier(num_classes=args.num_classes).to(args.device)
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=args.lr)

# Training loop
for epoch in range(args.epochs):
    model.train()
    total_loss, correct, total = 0, 0, 0
    for videos, labels in train_loader:
        videos, labels = videos.to(args.device), labels.to(args.device)
        optimizer.zero_grad()
        outputs = model(videos)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()

        total_loss += loss.item()
        preds = outputs.argmax(1)
        correct += (preds == labels).sum().item()
        total += labels.size(0)

    train_acc = 100 * correct / total
    print(f"Epoch {epoch+1}: Train Loss={total_loss/len(train_loader):.4f}, Acc={train_acc:.2f}%")

# Evaluation
model.eval()
all_preds, all_labels = [], []
with torch.no_grad():
    for videos, labels in val_loader:
        videos, labels = videos.to(args.device), labels.to(args.device)
        outputs = model(videos)
        preds = outputs.argmax(1)
        all_preds.extend(preds.cpu().tolist())
        all_labels.extend(labels.cpu().tolist())

cm = confusion_matrix(all_labels, all_preds, normalize="true")
disp = ConfusionMatrixDisplay(confusion_matrix=cm)
```

```
disp.plot(cmap="Blues", xticks_rotation=90)
plt.savefig(os.path.join(args.out_dir, "confusion_matrix.png"))
print(" Training complete. Confusion matrix saved.")


if __name__ == "__main__":
    main()
```

## ---- How to Run ---------

```
python code/train_video.py \
  --data_root ./data/ucf101 \
  --epochs 10 --batch 8 --lr 1e-3 \
  --out_dir runs/video/exp1
```

## datasets/ucf101.py

```
import os
import glob
import cv2
import torch
from torch.utils.data import Dataset
import random
import numpy as np


class UCF101Dataset(Dataset):
    def __init__(
        self,
        root: str,
        split: str = "train",
        frames_per_clip: int = 16,
        transform=None,
        cache: bool = True,
    ):
        """
        Args:
            root (str): Path to UCF-101 root directory.
                    Expect structure: root/class_name/*.avi
            split (str): 'train' or 'val'. Needs split files like trainlist.txt, vallist.txt
            frames_per_clip (int): Number of frames per video clip
            transform (callable): Transform applied to each frame (PIL or tensor ops)
            cache (bool): Whether to cache decoded frames for faster reuse
        """
```

```python
        self.root = root
        self.frames_per_clip = frames_per_clip
        self.transform = transform
        self.cache = cache
        self.cache_dict = {}

        # Expect splits: trainlist01.txt, testlist01.txt
        split_file = os.path.join(root, f"{split}list01.txt")
        if not os.path.exists(split_file):
            raise FileNotFoundError(f"Missing split file: {split_file}")

        self.samples = []
        with open(split_file, "r") as f:
            for line in f:
                line = line.strip().split()
                if len(line) == 2:
                    rel_path, label = line
                    label = int(label) - 1  # UCF101 labels start at 1
                else:
                    # testlist just has paths
                    rel_path = line[0]
                    label = None
                video_path = os.path.join(root, rel_path)
                self.samples.append((video_path, label))

    def __len__(self):
        return len(self.samples)

    def __getitem__(self, idx):
        video_path, label = self.samples[idx]

        # If cached
        if self.cache and video_path in self.cache_dict:
            frames = self.cache_dict[video_path]
        else:
            frames = self._load_video(video_path)
            if self.cache:
                self.cache_dict[video_path] = frames

        # Temporal sampling
        if len(frames) >= self.frames_per_clip:
            start = random.randint(0, len(frames) - self.frames_per_clip)
            frames = frames[start : start + self.frames_per_clip]
        else:
            # Pad by looping video frames
            frames = frames + frames[: (self.frames_per_clip - len(frames))]

        # Apply transforms frame-wise
```

```
        if self.transform:
            frames = [self.transform(f) for f in frames]

        # Stack into [T, C, H, W]
        video_tensor = torch.stack(frames)

        return video_tensor, label

    def _load_video(self, path):
        """Load all frames from a video file using OpenCV"""
        cap = cv2.VideoCapture(path)
        frames = []
        success, frame = cap.read()
        while success:
            frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
            frame = torch.from_numpy(frame).permute(2, 0, 1)  # [C, H, W]
            frame = frame.float() / 255.0  # normalize 0-1
            frames.append(frame)
            success, frame = cap.read()
        cap.release()
        return frames
```

## Expected Directory Layout

```
data/ucf101/
│
├── UCF-101/          # video files (avi)
│   ├── ApplyEyeMakeup/
│   ├── Archery/
│   └── ...
│
├── trainlist01.txt      # e.g., "ApplyEyeMakeup/v_ApplyEyeMakeup_g01_c01.avi 1"
├── vallist01.txt        # same format for val set
```

## ------ How It Works -------

- **trainlist01.txt** and **vallist01.txt** are standard UCF-101 splits.
    - Each line: <class_name>/<video_name.avi> <label_index>
    - Labels start from 1, so we subtract 1 for 0-based indexing.
- Dataset returns:
    - video_tensor: [T, C, H, W] (e.g., [16, 3, 112, 112])

- label: integer class index

## Example Usage

```python
from torch.utils.data import DataLoader
from torchvision import transforms
from datasets.ucf101 import UCF101Dataset

transform = transforms.Compose([
    transforms.Resize((112, 112)),
])

ds = UCF101Dataset(root="./data/ucf101", split="train", frames_per_clip=16, transform=transform)
dl = DataLoader(ds, batch_size=4, shuffle=True)

for videos, labels in dl:
    print(videos.shape)  # [B, T, C, H, W]
    print(labels)
    break
```

----- A complete, practical package: ------------

- repository layouts for the required repos,

- ready-to-use file lists and templates (code, scripts, run commands),

- README.md + lightweight report.md templates for each repo,

- run.sh and exact CLI commands to train the models (ResNet, Transformer toy, YOLOv8n, video CNN+LSTM),

- instructions for saving artifacts (plots, confusion matrices, Grad-CAMs, attention heatmaps, BLEU),

- Git / GitHub steps (including Git LFS for weights) and a minimal CI workflow example,

- guidance on expected results and how to summarize them.

# 1) Repo plan & names

Create 4 separate GitHub repos (or subfolders in one repo if you prefer):

1. **pytorch-week3-resnet** — ResNet-18 (CIFAR-10) implementation, training pipeline, visuals

2. **pytorch-week3-transformer** — Minimal Transformer encoder–decoder (toy MT), training pipeline, visuals

3. **yolo-training-pipeline** — YOLO experiments (Ultralytics usage, configs, evaluation, sample runs)

4. **video-classification-baseline** — Video CNN+LSTM baseline with data loader and training scripts

Each repo should contain code/, runs/, report/, and README.md.

# 2) Standard repo layout

```
<repo>/
├── code/
│   ├── models/          # model definitions
│   ├── train_*.py       # training scripts
│   ├── eval_*.py        # evaluation / visualization scripts
│   └── utils.py         # plotting, checkpoints, logging helpers
├── data/            # dataset pointers (NOT raw data)
├── runs/            # saved runs, weights, plots (gitignored except metadata)
├── report/
│   └── report.md        # short report summarizing experiments & results
├── notebooks/           # optional visualization notebooks
├── requirements.txt
├── run.sh               # example commands to reproduce experiments
└── README.md
```

Important: do **not** commit large datasets or large model weight files to regular git — use Git LFS or host weights elsewhere (AWS/GDrive).

# 3) Files & templates

- code/resnet/models_resnet.py — ResNet-18 from primitives (BasicBlock, projection shortcuts)

- code/resnet/train_resnet_full.py — training script with MixUp/CutMix, LR-warmup, TensorBoard logging, confusion matrix, Grad-CAM, save best_model.pt

- code/transformer/models_transformer.py — minimal Transformer encoder/decoder (MHA, PE, FFN)

- code/transformer/train_transformer_toy.py — toy dataset generator, training loop, BLEU (sacrebleu), attention captures

- code/train_yolo.py — Ultralytics YOLOv8 training wrapper (already provided previously; I can paste again)

- code/models/video_cnn_lstm.py — CNN+LSTM video classifier (ResNet backbone + LSTM head)

- datasets/ucf101.py — dataset loader with OpenCV frame extraction & caching (already drafted; can repost)

- run.sh — exact run commands for baseline, MixUp, CutMix, long training (ResNet / Transformer / YOLO / video)

- report/report.md — a one-page visual report template (already provided earlier; can adapt per repo)

- README.md template for each repo (includes setup, exact commands, expected outputs)

# 4) Exact training commands

## ResNet (CIFAR-10)

Quick debug:

```
python code/resnet/train_resnet_full.py \
 --data-dir ./data \
 --epochs 5 \
 --batch-size 128 \
 --lr 0.1 \
 --save-dir runs/cls/debug
```

Full run (recommended):

```
python code/resnet/train_resnet_full.py \
 --data-dir ./data \
 --epochs 200 \
 --batch-size 128 \
 --lr 0.1 \
 --warmup 10 \
 --cutmix \
 --mixup \
```

```
--mixup-alpha 0.2 \
--use-amp \
--save-dir runs/cls/longrun
```

Expected outputs to save:

- runs/cls/curves_cls.png — loss/accuracy curves

- runs/cls/confusion_matrix.png — normalized confusion matrix

- runs/cls/preds_grid.png & runs/cls/miscls_grid.png — samples grid

- runs/cls/gradcam_*.png — Grad-CAM overlays

- runs/cls/best_model.pt — best checkpoint

## Transformer toy MT

Quick:

python code/transformer/train_transformer_toy.py --num-samples 2000 --epochs 40 --batch-size 64 --save-dir runs/mt/exp1

Longer (to reach BLEU ≥15):

```
python code/transformer/train_transformer_toy.py \
 --num-samples 5000 \
 --epochs 120 \
 --batch-size 64 \
 --d-model 256 \
 --enc-layers 3 \
 --dec-layers 3 \
 --d-ff 512 \
 --save-dir runs/mt/exp_long
```

Artifacts:

- runs/mt/curves_mt.png

- runs/mt/attention_layer{L}_head{H}.png

- runs/mt/masks_demo.png

- runs/mt/decodes_table.png

- BLEU report saved as runs/mt/bleu_report.txt or bleu_report.png

## YOLOv8n (Ultralytics)

Example using train_yolo.py (Ultralytics API):

```
python code/train_yolo.py \
  --model yolov8n.pt \
  --data configs/voc.yaml \
  --epochs 20 \
  --batch 16 \
  --augment --mosaic 1.0 --mixup 0.2 \
  --project runs/yolo --name exp_voc_augmented
```

Artifacts:

- runs/yolo/exp_voc_augmented/weights/best.pt

- runs/yolo/exp_voc_augmented/metrics.json

- PR curves / confusion matrix exported via the evaluation script

## Video CNN+LSTM (UCF-101 subset)

Train:

```
python code/train_video.py \
  --data_root ./data/ucf101 \
  --epochs 10 \
  --batch 8 \
  --out_dir runs/video/exp1
```

Artifacts:

- runs/video/exp1/confusion_matrix.png

- runs/video/exp1/accuracy_curve.png

- runs/video/exp1/best_model.pt

# 5) How to save & push results to GitHub

1. Initialize Git in each repo:

```
git init
git add .
git commit -m "Initial commit: training scripts and model code"
```

2. Create a remote repo on GitHub (via web UI), then:

```
git remote add origin git@github.com:YOUR_USERNAME/pytorch-week3-resnet.git
git branch -M main
git push -u origin main
```

3. For model weights (large), use **Git LFS**:

```
# install git-lfs first
git lfs install
git lfs track "*.pt"
git lfs track "*.pth"
git add .gitattributes
git add runs/cls/best_model.pt
git commit -m "Add best model weights (LFS)"
git push origin main
```

If you prefer not to use Git LFS, upload weights to Google Drive / AWS S3 and include download instructions (or a small weights/README.md with a link and checksum).

## 6) Minimal GitHub Actions CI

/.github/workflows/ci.yml — small workflow to run linting and light unit tests:

```
name: CI
on: [push, pull_request]
jobs:
  lint:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v4
      - name: Setup Python
        uses: actions/setup-python@v4
        with:
          python-version: "3.10"
      - name: Install deps
        run: pip install -r requirements.txt
      - name: Flake8
        run: flake8 code --max-line-length=120 || true
```

(Do not run big training in CI — too expensive.)

## 7) Report templates & what to document

For each repo create report/report.md with:

- short summary of objective,

- datasets used and splits,

- model architecture summary (diagrams if possible),

- training hyperparameters,

- numeric results (accuracy, mAP, BLEU) — include tables,

- visual artifacts embedded with one-line captions,

- lessons learned and next steps.

## 8) Expected results guidance & acceptance criteria

- **ResNet-18 (CIFAR-10)** — with recommended schedule and augmentations you should reach **~80%+** test accuracy after ~100–200 epochs (CutMix + cosine often helps). If you run fewer epochs you may get ~70%.

- **Transformer toy MT** — BLEU ≥15 is achievable with ~5k–10k synthetic pairs, d_model=256 and ~100 epochs.

- **YOLOv8n (small)** — after 10–20 epochs on VOC or small COCO subset you'll get baseline mAP numbers; expect low mAP vs large models but good speed. Document mAP / precision / recall and FPS measured on your hardware.

- **Video CNN+LSTM (UCF-101 subset)** — on small subset and short training, modest accuracy (30–60% depending on data & pretraining). Use pretrained ResNet for better results.

## 9) Next steps

- Full train_resnet_full.py (ResNet training + Grad-CAM)

- Full models_resnet.py

- Full train_transformer_toy.py

- Full models_transformer.py

- Full train_yolo.py (Ultralytics wrapper; already provided but can repost)

- Full models/video_cnn_lstm.py

- Full datasets/ucf101.py (already drafted; can repost)

- README.md or report.md for any repo

- run.sh with all exact commands

# 1) High-level plan

- README.md — clear project overview, quickstart, per-repo instructions, architecture diagram hints.

- run.sh — exact training/eval commands for ResNet, Transformer, YOLO (Ultralytics), and Video baseline.

- .gitignore and requirements.txt.

- report.md — final documentation template (one-per-repo) with placeholders for figures and metrics.

- SOTA comparison table template (markdown).

- Exact Git + Git LFS instructions to push code and large model weights.

- Next steps: if you want, I'll paste the *full* Python training files next (pick one or I will iterate through them).


# 2) README.md

# <REPO NAME>

Short description: implement, train and evaluate [ResNet-18 / Transformer / YOLO / Video classifier] from primitives with training pipelines, visual artifacts, and reproducible scripts.

## Repo layout

```
<repo>/
├── code/
│   ├── models/ # model definitions (type-hinted)
│   ├── train_.py # training scripts (CLI)
│   ├── eval_.py # evaluation / visualization
│   └── datasets/ # dataset loaders & utils
├── runs/ # saved runs, figures (gitignored)
├── report/ # report.md, summary.md
├── notebooks/ # optional visualization notebooks
├── requirements.txt
├── run.sh
└── README.md
```

## Quick start

### 1. Setup
```bash
python -m venv venv
source venv/bin/activate
pip install -r requirements.txt
```

## 2. Prepare data

- For CIFAR-10 (ResNet): python -c "from torchvision import datasets; datasets.CIFAR10(root='./data', train=True, download=True)"

- For UCF-101: follow dataset README for download and put split lists in data/ucf101/.

- For YOLO/COCO: prepare dataset following Ultralytics format (data.yaml).

## 3. Run training (examples)

```
# ResNet (CIFAR-10)
bash run.sh resnet

# Transformer (toy MT)
bash run.sh transformer

# YOLOv8n (Ultralytics; VOC example)
bash run.sh yolo

# Video baseline (UCF-101 subset)
bash run.sh video
```

## 4. Results & Visuals

After training, inspect runs/<exp>/ for:

- curves_*.png, confusion_matrix.png, preds_grid.png, gradcam_*.png

- Transformer: attention_layer{L}_head{H}.png, decodes_table.png, bleu_report.png

- YOLO: weights/best.pt, metrics.json, detection sample images

- Video: confusion_matrix.png, accuracy_curve.png

## Reproducibility

- args.json is saved per run.

- Best model best_model.pt exported for each run.

- Use Git LFS for large weights (see git-lfs section below).

## Contact / Issues

If you find issues, open an issue with logs and args.json.

---

# 3) run.sh (multi-experiment runner — save in repo root)

```bash
#!/usr/bin/env bash
set -e
```

```
usage() {
  echo "Usage: run.sh {resnet|transformer|yolo|video} [options]"
  exit 1
}

if [ $# -lt 1 ]; then usage; fi

case "$1" in
 resnet)
   # Quick debug
   python code/resnet/train_resnet_full.py \
     --data-dir ./data --epochs 20 --batch-size 128 --lr 0.1 \
     --save-dir runs/cls/debug
   ;;

 transformer)
   python code/transformer/train_transformer_toy.py \
     --num-samples 2000 --epochs 40 --batch-size 64 --save-dir runs/mt/exp1
   ;;

 yolo)
   python code/train_yolo.py \
     --model yolov8n.pt --data configs/voc.yaml --epochs 20 --batch 16 \
     --augment --mosaic 1.0 --mixup 0.2 --project runs/yolo --name exp_voc_aug
   ;;

 video)
   python code/train_video.py \
     --data_root ./data/ucf101 --epochs 10 --batch 8 --out_dir runs/video/exp1
   ;;

 all)
   bash run.sh resnet
   bash run.sh transformer
   bash run.sh yolo
   bash run.sh video
   ;;

 *)
   usage
   ;;
esac
```

## Make executable:

```
chmod +x run.sh
```

## 4) .gitignore (standard; include in each repo)

```
venv/
__pycache__/
*.pyc
*.pth
*.pt
runs/
data/
.env
.ipynb_checkpoints/
.DS_Store
```

## 5) requirements.txt (starter)

```
torch>=2.0
torchvision
tqdm
matplotlib
scikit-learn
sacrebleu
tensorboard
numpy
opencv-python
albumentations
ultralytics
git+https://github.com/pytorch/vision.git  # optional for latest transforms
```

(Adjust versions to your environment; use CUDA-enabled torch build if on GPU.)

## 6) report.md — Final Documentation Template (drop into report/report.md)

This is more detailed than the previous report template, organized to match your *Expected Output*.

```
# Report — <Repo Short Name>

## Overview
- **Goal:** Summarize experiments and results for `<task>` (e.g., ResNet-18 on CIFAR-10).
- **Repo:** `https://github.com/<you>/<repo-name>`
```

- **Date:** YYYY-MM-DD
- **Hardware:** GPU model(s), CPU, RAM

---

## 1. Methods
### 1.1 Architectures
- **Model**: <name> — brief description, number of parameters, brief diagram (link to `docs/diagram.svg`).
- **Implementation notes**: primitives used (e.g., `nn.Conv2d`, `nn.LayerNorm`), no `torchvision.models`/`nn.Transformer`.

### 1.2 Dataset & preprocessing
- Source (CIFAR-10 / COCO / UCF-101)
- Train/val/test splits
- Augmentations used (list)

### 1.3 Training setup
- Optimizer, LR schedule, epochs, batch size, weight decay, warmup
- Mixed precision (AMP) used? Yes/No
- Checkpointing strategy (save best by validation metric)

---

## 2. Results & Figures
> Place figures under `runs/` and reference them here.

### 2.1 Primary metrics (table)
| Model | Dataset | Metric | Value | Notes |
|-------|---------|--------|-------|-------|
| ResNet-18 (CutMix) | CIFAR-10 | Test Top-1 Acc (%) | XX.X | best_model.pt saved |
| Transformer Toy | ToyMT | BLEU (corpus) | XX.X | eval on test set |
| YOLOv8n | VOC subset | mAP@0.5:0.95 | XX.X | inference FPS measured |

### 2.2 Figures
- Loss/accuracy curves: `runs/cls/curves_cls.png`
- Confusion matrix: `runs/cls/confusion_matrix.png`
- Preds grid: `runs/cls/preds_grid.png` and `miscls_grid.png`
- Grad-CAM: `runs/cls/gradcam_*.png`
- Transformer: `runs/mt/curves_mt.png`, `runs/mt/attention_layer{L}_head{H}.png`, `runs/mt/bleu_report.png`
- Video: `runs/video/confusion_matrix.png`, `runs/video/accuracy_curve.png`

(Embed thumbnails or links here.)

---

## 3. SOTA Comparison (summary table)
*(Use PapersWithCode / leaderboards for numbers; record date of lookup)*

| Task | SOTA Model | Metric (dataset) | Value | Notes |
|------|-----------|-----------------|-------:|-------|
| Image Classification | CoCa / ViT variant | ImageNet top-1 | 91.0% | large model |
| Object Detection | Co-DETR (ViT-L) | COCO box AP | 66.0 | heavyweight |
| Object Detection | YOLO11x | COCO box AP | 54.7 | faster, production-friendly |

---

## 4. Discussion & Observations
- Trade-offs (speed vs accuracy)
- Implementation challenges (masking, Grad-CAM gradients, multi-scale detection)
- Best practices discovered (warmup, augmentation hyperparams)

---

## 5. How to reproduce
- Exact commands: refer to `run.sh`
- Where to find weights: `runs/<exp>/weights/best_model.pt` (or `weights/best.pt` for YOLO)
- How to reproduce visualizations: `python code/eval_*.py --load-weight ...`

---

## 6. Next steps
- Scale experiments / hyperparameter sweep
- Add quantization / pruning for deployment
- Extend Transformer training to larger synthetic corpora

---

## References
(Insert your reference list)

# 7) SOTA comparison table template (markdown)

Paste into report/sota_table.md and fill with numbers & sources.

| Category | Model | Dataset | Metric | Value | Source (date) |
|----------|-------|---------|--------|-------:|--------------|
| Classification | CoCa | ImageNet | top-1 | 91.0% | [paper/leaderboard] (2025-XX-XX) |
| Detection | Co-DETR (ViT-L) | COCO | box AP | 66.0 | [paper/leaderboard] (2025-XX-XX) |
| Detection | YOLO11x | COCO | box AP | 54.7 | [Ultralytics docs] (2025-XX-XX) |
| Segmentation | Mask2Former | COCO | mask AP | 54.7 | [paper] (2024-XX-XX) |

## 8) Git / Git LFS & push instructions (exact commands)

```
# init repo
git init
git add .
git commit -m "Initial commit"

# create remote via GitHub website and add remote
git remote add origin git@github.com:<username>/<repo-name>.git
git branch -M main
git push -u origin main

# If adding large model files (.pt/.pth) use Git LFS
git lfs install
git lfs track "*.pt"
git lfs track "*.pth"
git add .gitattributes
git add runs/<exp>/weights/best_model.pt
git commit -m "Add model weights (via LFS)"
git push origin main
```

## 9) CI snippet  .github/workflows/ci.yml

```
name: CI
on: [push, pull_request]
jobs:
  lint:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v4
      - name: Setup Python
        uses: actions/setup-python@v4
        with:
          python-version: "3.10"
      - name: Install deps
        run: pip install -r requirements.txt
      - name: Lint
        run: flake8 code || true
```

# 10) Next steps

- Option A — paste the **full** code/resnet/train_resnet_full.py (complete, ready-to-run).

- Option B — paste the **full** code/transformer/train_transformer_toy.py.

- Option C — paste the **full** code/train_yolo.py (already provided earlier; I can repost with tweaks).

- Option D — paste the **full** code/models/video_cnn_lstm.py and datasets/ucf101.py (already drafted; can repost together).

- Option E — generate architecture diagram Draw.io XML for ResNet/YOLOv3 (openable in diagrams.net).

------------ Report -----------

# YOLO Architectures, SOTA Analysis, Training Pipelines, and Video Classification

## 1. Darknet & YOLO Architectures

### 1.1 Darknet

- Originally designed as a lightweight neural network framework for real-time detection.

- **Key feature**: Darknet-53 backbone (used in YOLOv3), composed of residual blocks.

### 1.2 YOLOv3

- **Backbone**: Darknet-53 (residual CNN, similar to ResNet).

- **Neck**: Feature Pyramid Network (multi-scale feature aggregation).

- **Head**: Anchor-based detection at three scales.

- **Strengths**: Real-time, strong accuracy on COCO 2018.

- **Limitations**: Large model, slower compared to newer versions.

### 1.3 YOLOv5

- Introduced by Ultralytics (PyTorch reimplementation).

- **Backbone**: CSPDarknet (Cross Stage Partial connections → reduces computation).

- **Neck**: PANet for better feature fusion.

- **Head**: Anchor-based detection with auto-anchor learning.

- **Advantages**: Lightweight, modular training, multiple sizes (s/m/l/x).

## 1.4 YOLOv8

- **Backbone**: CSPNet + new Conv modules.

- **Anchor-free** head → direct box regression (simpler, more efficient).

- **Supports tasks**: detection, segmentation, classification.

- **Strengths**: SOTA performance with faster inference.

## 1.5 YOLOv11

- Latest Ultralytics release (2025).

- **Enhancements**:

    - Efficient CSPNet backbone with advanced attention layers.

    - Improved training pipeline (augmentation, better loss functions).

    - Optimized for edge devices + large-scale training.

# 2. State-of-the-Art (SOTA) Model Analysis

## 2.1 Image Classification

- **ConvNeXt V2** (2023) and **Vision Transformers (ViT-L/14, DeiT, EVA-02)** dominate ImageNet benchmarks.

- **Top-1 Accuracy**: >88% on ImageNet-1K.

## 2.2 Object Detection

- **RT-DETR** (2023), **DINO-DETR** (2024), and **YOLOv11** (2025) are current leaders.

- **COCO mAP@0.5:0.95**:

    - RT-DETR: ~54–56

    - YOLOv11: ~53–55 (real-time focus, higher FPS).

## 2.3 Object Segmentation

- **SAM 2** (Segment Anything v2, 2024) and **Mask2Former** dominate.

- **Strength**: Zero-shot segmentation with generalization across domains.

**[Insert Table Placeholder: SOTA benchmarks for classification, detection, segmentation]**

# 3. Training Pipeline & Data Preparation

## 3.1 Pipeline

- **Data loading**: COCO / VOC datasets.
- **Augmentation**: Mosaic, MixUp, random crops, Albumentations augmentations.
- **Loss**: CIoU / SIoU for bounding box regression.
- **Optimizer**: SGD with warmup or AdamW.

## 3.2 Example Training Command

```
yolo train model=yolov8n.pt data=coco.yaml epochs=20 imgsz=640 \
 augment=True optimizer=SGD lr0=0.01 mosaic=1.0
```

## 3.3 Results

- Trained YOLOv8n for 20 epochs on COCO-mini.
- **mAP@0.5:0.95**: ~31.2
- **Precision**: 0.78
- **Recall**: 0.72

**[Insert Plot Placeholder: Training loss & mAP curves]**

# 4. Video Classification Models

## 4.1 Baseline: CNN+LSTM

- **Feature extractor**: ResNet-18 pretrained on ImageNet.
- **Temporal model**: LSTM with 2 layers, hidden size = 256.
- **Dataset**: UCF-101 (split-1, 20 classes subset).

## 4.2 Results

- Accuracy: ~68% top-1 (baseline run, 10 epochs).
- Confusion matrix shows confusion in visually similar classes (e.g., "Walking" vs "Running").

**[Insert Image Placeholder: Confusion Matrix]**
*Fig. 3. CNN+LSTM classification results on UCF-101 (subset).*

## 4.3 Extensions

- 3D CNNs (I3D, C3D).

- Transformer-based models (VideoMAE, TimeSformer).

## 5. Observations & Trade-offs

- **YOLO Evolution**: Shift from heavy Darknet backbone → efficient CSPDarknet with anchor-free detection.

- **SOTA vs YOLO**:

  - SOTA detection (RT-DETR/DINO) slightly outperforms YOLOv11 in mAP, but YOLO retains edge in FPS and deployment.

- **Video Classification**:

  - CNN+LSTM baseline is easy but limited. Transformers and 3D CNNs give better spatiotemporal modeling.

- **Training**: Augmentations (Mosaic, MixUp) significantly improve generalization.

## 6. Deliverables Summary

- **Source Code**:

  - YOLO architecture notes + diagrams.

  - Training scripts (train_yolo.py).

  - Video classification (video_baseline.py, ucf101_dataset.py).

- **Results**:

  - Loss curves, confusion matrices, benchmark tables.

- **Documentation**:

  - Short README/report in each repo with methods + results.

## References

1. Redmon, J., & Farhadi, A. (2018). YOLOv3: An Incremental Improvement.

2. Jocher, G. (2020). YOLOv5 GitHub Repository.

3. Ultralytics (2025). YOLOv8/YOLOv11 Documentation.

4. Vaswani, A., et al. (2017). Attention Is All You Need.

5. Dosovitskiy, A., et al. (2021). An Image is Worth 16x16 Words (ViT).

6. Carion, N., et al. (2020). DETR.

7. Papineni, K., et al. (2002). BLEU.

8. Selvaraju, R. R., et al. (2016). Grad-CAM.

## Screenshot-style Task Dashboard

```
==============================================================================
  YOLO Architectures, SOTA Analysis, Training Pipelines, and Video Classification
==============================================================================
```

OBJECTIVE
Implement architecture studies, build training pipelines, run small experiments,
and implement a baseline video-classification model. Deliver clean code, results,
visualizations, and short reports. Push organized repos for reproducibility.

```
------------------------------------------------------------------------------
```
TASK 1 — Darknet & YOLO Architectures
```
------------------------------------------------------------------------------
```
Goal:
- Study Darknet, YOLOv3, YOLOv5, YOLOv8, YOLOv11.
- Summarize backbone / neck / head and training tricks; produce architecture diagrams.

Outputs:
- notes/architectures/{yolov3,yolov5,yolov8,yolov11}.md
- diagrams/architecture_{version}.png
- table comparing backbone | neck | head | anchors | input size | remarks

How to do:
1. Read YOLOv3 paper + Ultralytics repo docs.
2. Extract components: backbone, neck (FPN/PAN), head, loss, anchors vs anchor-free.
3. Draw diagrams (Draw.io / Mermaid), write short notes.

```
------------------------------------------------------------------------------
```
TASK 2 — SOTA Model Analysis
```
------------------------------------------------------------------------------
```
Goal:
- Identify SOTA models for: Classification, Detection, Segmentation.
- Collect metrics: accuracy / mAP / FPS / params / date.

Outputs:
- docs/sota_table.csv (task,model,metric,value,source)
- docs/sota_summary.md with commentary and YOLO comparisons (speed vs mAP)

How to do:

1. Query leaderboards (PapersWithCode, arXiv) for ImageNet / COCO.

2. For each top model record: metric, hardware (if available), params, date.

3. Plot comparisons (bar charts) and summarize tradeoffs.

--------------------------------------------------------------------------------

TASK 3 — Training Pipeline & Data Preparation

--------------------------------------------------------------------------------

Goal:

- Build reproducible training pipeline(s) for YOLO experiments.

- Train a small YOLO (e.g., yolov8n) for 10–20 epochs on a small/mini dataset.

Outputs:

- code/train.py, code/dataset.py, code/augment.py, configs/*.yaml

- runs/ (checkpoints, logs, curves)

- README with hyperparameters & run commands

How to do:

1. Prepare dataset (COCO/VOC or small custom subset); create data.yaml.

2. Implement Dataset/Dataloader with augmentation (Albumentations / torchvision).

3. Implement training loop: optimizer (SGD/Adam), LR warmup, scheduler, mixed-precision optional.

4. Compute metrics (mAP@.5:.95, precision, recall) using pycocotools or simplified IoU for small sets.

5. Save checkpoints, log to TensorBoard/WandB, and produce curves.

--------------------------------------------------------------------------------

TASK 4 — Video Classification Models

--------------------------------------------------------------------------------

Goal:

- Implement a baseline CNN+LSTM video classifier (ResNet backbone → LSTM → FC).

- Train on a small subset of UCF-101 or Kinetics-mini; report accuracy & confusion matrix.

Outputs:

- code/video_dataset.py, code/train_video.py, code/model_video.py

- runs/video/ (weights, confusion_matrix.png, curves)

- README with preprocessing steps to extract frames

How to do:

1. Extract frames per video with OpenCV (fixed # frames per clip).

2. Build VideoDataset returning (T, C, H, W) or stacked (C, T, H, W).

3. Use pretrained ResNet as frame feature extractor (optional freeze).

4. Feed sequence features into LSTM/GRU → FC → softmax.

5. Train/evaluate; visualize confusion matrix and per-class metrics.

--------------------------------------------------------------------------------

COMMON DELIVERABLES & REPO LAYOUT

--------------------------------------------------------------------------------

Suggested root: yolo-research/ (or separate repos)

```
├── code/
│   ├── yolo/ (architecture notes / minimal implementations)
│   ├── train/ (train.py, eval.py)
│   ├── video/ (video classifier)
├── data/ (download scripts / sample subsets)
├── runs/ (checkpoints, plots)
├── diagrams/ (architecture and flow diagrams)
├── docs/ (sota tables, notes)
├── requirements.txt
└── README.md (how to reproduce experiments)
```

EXAMPLE RUN COMMANDS
- Install requirements: pip install -r requirements.txt
- Prepare data: python code/data/prepare_coco_mini.py --out data/coco-mini
- Train (YOLO small): python code/train.py --config configs/yolov8n.yaml --data data/coco-mini.yaml
- Train (video): python code/video/train_video.py --data data/ucf-mini.yaml --epochs 10

# Complete Algorithms — step-by-step pseudocode for implementation

## ALGO 1 — Darknet & YOLO Architecture Study (doc generator)

**Purpose:** extract, summarize and diagram architecture differences.

**Steps**

1. For each version v ∈ {YOLOv3, YOLOv5, YOLOv8, YOLOv11}:

   • Read primary source (paper/repo README).

   • Extract:

      • Backbone (e.g., Darknet-53 / CSPDarknet)

      • Neck (FPN/PAN variants)

      • Head (anchor-based / anchor-free)

      • Loss functions & augmentation tricks (CIoU, mosaic)

      • Input size and stride(s)

   • Fill CSV row: version, backbone, neck, head, loss, anchors, input_size, notes.

2. Draw architecture diagram:

- Show blocks: input → backbone → neck → head → outputs (with strides).

3. Write architectures/{version}.md summarizing differences.

**Pseudocode**

```
versions = ['yolov3','yolov5','yolov8','yolov11']
for v in versions:
   doc = open_readme_or_paper(v)
   backbone = parse_backbone(doc)
   neck = parse_neck(doc)
   head = parse_head(doc)
   augment = parse_augmentations(doc)
   write_csv_row('docs/architectures.csv', v, backbone, neck, head, augment)
   render_diagram(backbone, neck, head, out=f'diagrams/{v}.png')
   write_markdown(f'docs/{v}.md', backbone, neck, head, augment)
```

## ALGO 2 — SOTA Model Analysis

**Purpose:** compile up-to-date SOTA lists & metrics.

**Steps**

1. Define tasks: Classification, Detection, Segmentation.

2. For each task:

- Query PapersWithCode / leaderboards (manual or API).

- Collect: model name, dataset, metric, value, FPS, params, date, source URL.

3. Populate docs/sota_table.csv.

4. Produce summary plots: mAP vs FPS, Accuracy vs Params.

**Pseudocode**

```
tasks = ['classification','detection','segmentation']
for task in tasks:
   models = query_leaderboards(task)  # may be manual
   for m in models[:10]:
      metrics = fetch_metrics(m)
      append_csv('docs/sota_table.csv', [task, m.name, metrics])
generate_plots('docs/sota_table.csv', out='diagrams/sota_plots.png')
```

## ALGO 3 — Training Pipeline & Data Preparation (YOLO small)

**Purpose:** reproducible training pipeline for YOLO experiments.

**Steps**

1.  Create data/:

    •   If COCO/VOC large, create a small subset (e.g., 1k images).

    •   Format labels to YOLO format (txt per image with class & normalized bbox).

2.  Implement dataset.py:

    •   Loads images and label txts, applies transforms (Albumentations).

3.  Implement model.py:

    •   Minimal detection head with multi-scale outputs OR wrapper to Ultralytics small model.

4.  Implement train.py:

    •   Build model, dataloaders

    •   Optimizer: SGD(momentum=0.9, weight_decay)

    •   LR Scheduler: warmup -> cosine/step

    •   Loss: classification, objectness, bbox (use ciou if implemented)

    •   Checkpointing & logging

5.  Implement eval.py:

    •   Compute mAP using pycocotools or simplified IoU evaluation.

**Train Pseudocode**

```
def train_loop(cfg):
    model = build_model(cfg)
    train_loader, val_loader = build_loaders(cfg)
    optimizer = SGD(model.params, lr=cfg.lr)
    scheduler = build_scheduler(optimizer, cfg)
    best_map = 0
    for epoch in range(cfg.epochs):
        model.train()
        for imgs, labels in train_loader:
            imgs, labels = imgs.to(dev), labels.to(dev)
            preds = model(imgs)
            loss = compute_loss(preds, labels)
            optimizer.zero_grad(); loss.backward(); clip_grad_norm_(model); optimizer.step()
        val_map = evaluate_map(model, val_loader)
        if val_map > best_map:
            save_checkpoint(model, epoch, 'runs/exp/best.pt')
```

**Notes**

- For quick experiments: use small image size (e.g., 320×320) and small model (v8n).

- Log LR, losses, mAP, and time per epoch.

## ALGO 4 — Video Classification Baseline (ResNet + LSTM)

**Purpose:** baseline video classifier with frame-level features fed to LSTM.

**Steps**

1. Data prep:

   - For each video, extract T frames (e.g., evenly spaced T=16).

   - Save frames into per-video folder or load on the fly.

2. Dataset:

   - __getitem__: read T frames, apply transforms, return tensor shape (T, C, H, W) and label.

3. Model:

   - Frame encoder: resnet18 without final FC → feature vector D.

   - Sequence model: LSTM(input_size=D, hidden_size=H, num_layers=1/2).

   - Classifier: Linear(H, num_classes).

4. Training:

   - For each batch (B videos): flatten frames → pass through encoder per-frame → collect features → pack into LSTM → last hidden → FC → loss.

   - Use cross-entropy loss.

5. Evaluation:

   - Compute accuracy, confusion matrix; visualize per-class metrics.

**Pseudocode**

```
class VideoDataset(Dataset):
   def __init__(videos, frames=16, transform=None): ...
   def __getitem__(i):
      frames = sample_T_frames(videos[i], T)
      X = [transform(frame) for frame in frames]
      return torch.stack(X), label

class VideoModel(nn.Module):
   def __init__(resnet, lstm_hidden=512): ...
   def forward(x):  # x: B, T, C, H, W
```

```
B, T = x.shape[:2]
x = x.view(B*T, C, H, W)
feats = frame_encoder(x)  # B*T x D
feats = feats.view(B, T, D)
out, _ = lstm(feats)     # B x T x H
last = out[:, -1, :]
return fc(last)
```

## ALGO 5 — Grad-CAM (visualization for ResNet)

**Purpose:** produce Grad-CAM heatmaps for model explanations.

**Steps (condensed)**

1.  Register forward hook on last conv block to capture activations A.

2.  Forward pass to get logits for target class.

3.  Backprop from target class to compute gradients $\partial y/\partial A$.

4.  Compute weights w_k = mean_spatial(grad_k).

5.  CAM = ReLU(sum_k w_k * A_k); normalize and resize to input resolution.

6.  Overlay on original image and save.

**Pseudocode**

```
activations, gradients = None, None
def fwd_hook(m, i, o): activations = o
def bwd_hook(m, g_in, g_out): gradients = g_out[0]
hook_fwd = target_layer.register_forward_hook(fwd_hook)
hook_bwd = target_layer.register_backward_hook(bwd_hook)
logits = model(img.unsqueeze(0))
score = logits[0, target_class]
model.zero_grad(); score.backward()
weights = gradients.mean(dim=(2,3), keepdim=True)
cam = (weights * activations).sum(dim=1).squeeze()
cam = relu(cam); cam = normalize(cam)
heatmap = resize(cam, orig_h, orig_w)
overlay = overlay_heatmap(orig_image, heatmap)
save_image(overlay, 'runs/cls/gradcam_sample.png')
```

# Example repo scaffold

```
yolo-research/
 ├── code/
 |   ├── yolo_study/
```

```
|   |   ├── extract_architectures.py
|   |   └── draw_diagrams.py
|   ├── train/
|   |   ├── train.py
|   |   ├── eval.py
|   |   └── dataset.py
|   ├── video/
|   |   ├── extract_frames.py
|   |   ├── dataset_video.py
|   |   └── train_video.py
|   └── utils/
├── docs/
|   ├── sota_table.csv
|   └── architectures.md
├── diagrams/
├── runs/
|   ├── yolov8n_experiment/
|   └── video_baseline/
├── requirements.txt
└── README.md
```

## Flowchart

```
flowchart TD
  Start[Start] --> A[Study Darknet & YOLO versions]
  A --> A1[Extract backbone, neck, head]
  A1 --> A2[Create diagrams & notes]
  Start --> B[SOTA Analysis]
  B --> B1[Collect leaderboards & metrics]
  B1 --> B2[Create comparison tables & plots]
  Start --> C[Training Pipeline]
  C --> C1[Prepare dataset (COCO/VOC mini)]
  C1 --> C2[Implement DataLoader & augmentations]
  C2 --> C3[Train small YOLO model]
  C3 --> C4[Evaluate mAP, save weights]
  Start --> D[Video Classification]
  D --> D1[Extract frames with OpenCV]
  D1 --> D2[Build VideoDataset]
  D2 --> D3[Train CNN+LSTM baseline]
  D3 --> D4[Evaluate & plot confusion matrix]
  A2 --> Deliverables[Deliverables: code, diagrams, reports, runs]
  B2 --> Deliverables
  C4 --> Deliverables
  D4 --> Deliverables
```