# PROJECT REPORT

Lecture: Video Compression Techniques

Lecturer: Kozyri Maria

Project Title:

Image and Video Compression using CompressAI

Contents:

- Environment Setup Description
- Part 1:
    - Model Benchmark Tests
    - Ablation Study
    - Evaluation
- Part 2(BONUS):
    - Bonus 1-CompressAI evaluation on Video Sequences
    - Bonus 2-Encoding-Decoding time calculation per Frame
    - Bonus 3-VMAF correlation

*Problems Encountered will be showcased in each section*

University of Thessaly

Academic Year 2025-2026

**Overview:** The project focuses on evaluating learned image and video compression methods using CompressAI and extending the study with BONUS experiments (Video-for-Machines, Complexity & Resources, and VMAF correlation). The main objectives were:

- Benchmark several learned image codecs on the Kodak dataset
- Implement an ablation study exploring the impact of preprocessing
- Extend to video sequences (UVG dataset) to measure performance for machine vision tasks (object detection accuracy)
- Collect computational and perceptual metrics for a complete evaluation

https://github.com/InterDigitalInc/CompressAI

https://r0k.us/graphics/kodak/

https://ultravideo.fi/dataset.html

## Environment Setup:

A virtual environment with Python version 3.11 was created in order to ensure stability between all package versions. Using pip the following packages were installed and finally exported to a "requirements.txt" file for reusability.

| Package Name | Version | Package Name | Version |
|---|---|---|---|
| aiohappyeyeballs | 2.6.1 | nvidia-nvjitlink-cu12 | 12.9.86 |
| aiohttp | 3.13.1 | nvidia-nvtx-cu12 | 12.1.105 |
| aiosignal | 1.4.0 | opencv-python | 4.9.0.80 |
| async-timeout | 5.0.1 | packaging | 25 |
| attrs | 25.4.0 | pandas | 2.3.3 |
| certifi | 2022.12.7 | pillow | 11.3.0 |
| charset-normalizer | 2.1.1 | propcache | 0.4.1 |
| compressai | 1.2.8 | psutil | 7.1.1 |
| contourpy | 1.3.2 | py-cpuinfo | 9.0.0 |
| cycler | 0.12.1 | pybind11 | 3.0.1 |
| einops | 0.8.1 | pyparsing | 3.2.5 |
| filelock | 3.19.1 | python-dateutil | 2.9.0.post0 |
| fonttools | 4.60.1 | pytorch-msssim | 1.0.0 |
| frozenlist | 1.8.0 | pytz | 2025.2 |
| fsspec | 2025.9.0 | PyYAML | 6.0.3 |
| idna | 3.4 | requests | 2.28.1 |
| Jinja2 | 3.1.6 | scipy | 1.15.3 |
| kiwisolver | 1.4.9 | seaborn | 0.13.2 |
| MarkupSafe | 2.1.5 | six | 1.17.0 |

| | | | |
|---|---|---|---|
| matplotlib | 3.10.7 | sympy | 1.14.0 |
| mpmath | 1.3.0 | tomli | 2.3.0 |
| multidict | 6.7.0 | torch | 2.2.0 |
| networkx | 3.4 | torch-geometric | 2.7.0 |
| numpy | 1.26.4 | torchvision | 0.17.0 |
| nvidia-cublas-cu12 | 12.1.3.1 | tqdm | 4.67.1 |
| nvidia-cuda-cupti-cu12 | 12.1.105 | triton | 2.2.0 |
| nvidia-cuda-nvrtc-cu12 | 12.1.105 | typing_extensions | 4.15.0 |
| nvidia-cuda-runtime-cu12 | 12.1.105 | tzdata | 2025.2 |
| nvidia-cudnn-cu12 | 8.9.2.26 | ultralytics | 8.3.30 |
| nvidia-cufft-cu12 | 11.0.2.54 | ultralytics-thop | 2.0.18 |
| nvidia-curand-cu12 | 10.3.2.106 | urllib3 | 1.26.13 |
| nvidia-cusolver-cu12 | 11.4.5.107 | utils | 1.0.2 |
| nvidia-cusparse-cu12 | 12.1.0.106 | xxhash | 3.6.0 |

_Problems Encountered: Several incompatibility issues arose due to numpy>=2.0, which was downgraded to maintain compatibility with compressai==1.2.8._

## Image Compression Baseline:

Implemented run_compressai.py, which evaluates multiple pretrained image codecs on the Kodak dataset:

- bmshj2018_factorized
- bmshj2018_hyperprior
- mbt2018_mean
- mbt2018
- cheng2020_anchor
- cheng2020_attn

A Quality ladder of 6 levels (Q=1,2,3,4,5,6) was set to all models. Computed standard metrics:

- BPP,
- PSNR

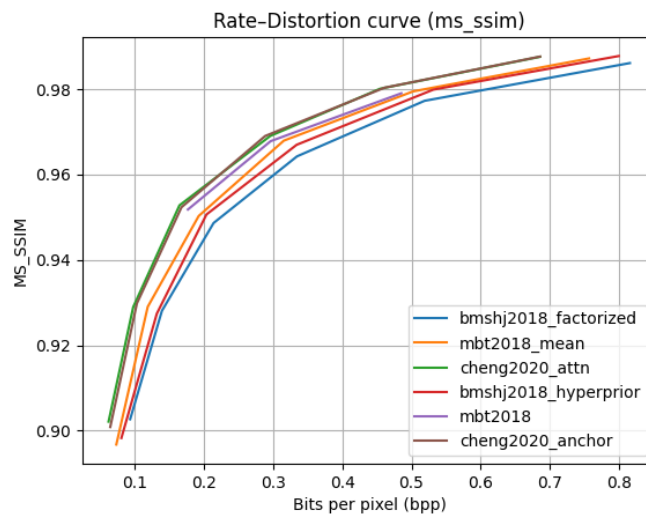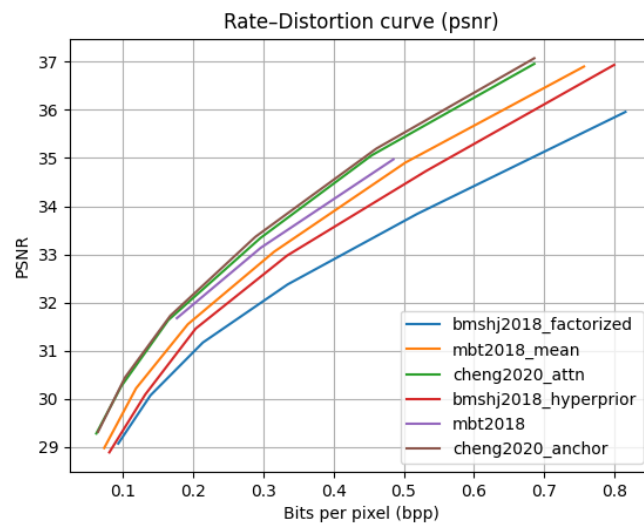- MS-SSIM
- encoding/decoding times
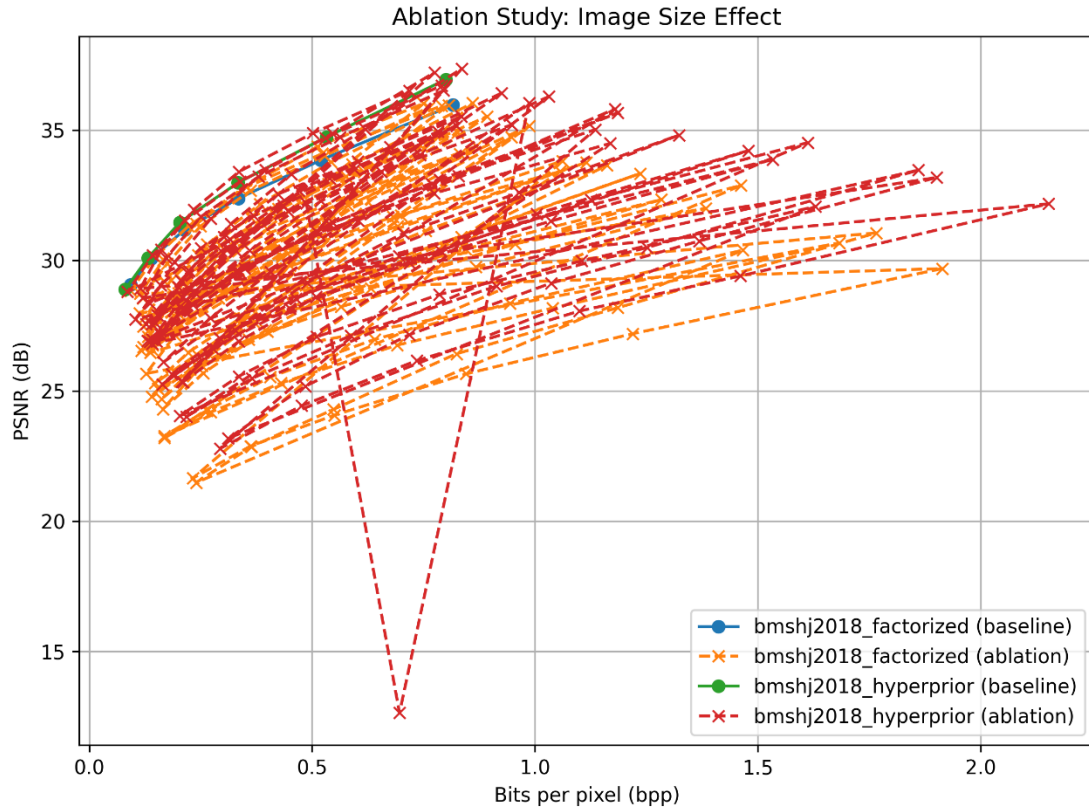
Exported results to "*results/image_rd_kodak.csv*"

As for plotting, "plot_rd_curves.py" generates:

- Rate–Distortion (RD) curves for PSNR and MS-SSIM
- Combined comparison between baseline and ablation runs

Outcome: Successfully reproduced RD curves consistent with the CompressAI reference results.

**Ablation Study:** Implemented ablation.py, resizing input images (e.g., to 192×192) to study the effect of resolution. Compared against the baseline via an extended plotting script:



Ablation Study: Image Size Effect

The main observation was that smaller inputs yield slightly lower PSNR at equal bitrates, confirming the expected trade-off between spatial detail and compressibility.

# BONUS-1: Video for Machines

Developed bonus.py to evaluate CompressAI on video sequences. Two datasets were used. Specifically, with qualities of 1080p and UVG clip names Bosphorus.yuv, HoneyBee.yuv, 100 frames each. The models that proved to be compatible with GPU (cuda) acceleration were bmshj2018_factorized and cheng2020_attn.

For each frame:

- Compressed/decompressed frame
- Computed PSNR, MS-SSIM, BPP
- Detected objects with YOLOv8n → number of detections used as proxy for mAP
- Generated Accuracy-vs-Rate curves

- Plot: plots/bonus1_accuracy_vs_rate.png X-axis: Bitrate (bpp)/Y-axis: Detection accuracy (mean detections per frame)

Outcome: Results show that detection performance degrades with aggressive compression; Cheng2020 maintains slightly higher accuracy at equivalent bitrate.

## BONUS-2: Complexity & Resources

The same script collected:

- Encoding and decoding times (per frame)
- CPU RAM and GPU VRAM usage (via psutil and torch.cuda.memory_allocated)
- Produced plot: plots/bonus2_rate_vs_time.png
- X-axis: Bitrate (bpp)
- Y-axis: Encoding time (s)

Observation: cheng2020_attn exhibits much higher computational complexity due to its autoregressive entropy coding. bmshj2018_factorized is 3–5× faster with smaller memory footprint.

## BONUS-3: VMAF Correlation

1. Integrated ffmpeg/libvmaf metric computation
2. Calculated VMAF for representative frames (every 10th frame) to correlate perceptual quality with PSNR
3. Generated scatter plot: plots/bonus3_vmaf_vs_psnr.png
4. X-axis: PSNR (dB)
5. Y-axis: VMAF (0–100)

Observation: Strong positive correlation ($R^2 \approx 0.9$). PSNR remains a reasonable indicator of perceived quality for high-quality regimes.

Problems Encountered and Solutions

| Issue | Description | Solution |
|---|---|---|
| FileNotFoundError: Zone.Identifier | Caused by Windows metadata streams in image folder | Filtered filenames using .endswith((".png",".jpg",".jpeg")) |

| | | |
|---|---|---|
| KeyError: 'strings' | Some CompressAI outputs differed between models | Fixed by aligning expected keys across versions and ensuring correct metric call |
| Invalid quality "9" | Models support Q∈[1,8] | Limited QUALITIES=[1–6] for consistency |
| Tensor type mismatch (CUDA vs CPU) | Occurred when image or model were on different devices | Ensured .to(device) applied consistently |
| Out of memory (CUDA error) | Full-resolution autoregressive models overloaded GPU | Moved CompressAI to CPU, added torch.cuda.empty_cache() |
| Autoregressive model slow execution | cheng2020_attn sequential entropy coder cannot parallelize | Accept slower runtime; added progress bars (tqdm) and checkpoint saving |
| VMAF computation empty results | Missing libvmaf model file or incompatible ffmpeg build | Installed ffmpeg with libvmaf and updated compute function |
| Very long runtime at 1080p | Each frame takes 30–60 s on CPU | Retained full resolution but introduced safe checkpointing |
| Dimension mismatch in video compression | Neural compression models require input dimensions divisible by 64 (1080p → 1088×1920) | Implemented padding functions pad_to_multiple() and unpad_tensor() with reflection padding |
| GPU utilization limitations | cheng2020_attn model's autoregressive entropy coder runs sequentially on CPU despite GPU availability | Switched to GPU-compatible models: bmshj2018_factorized, bmshj2018_hyperprior, mbt2018_mean |
| Memory fragmentation | Long-running experiments caused GPU memory fragmentation and gradual performance degradation | Added regular torch.cuda.empty_cache() calls and checkpoint-based memory management |
| YUV file reading bottleneck | Reading YUV420 frames required manual byte parsing and chroma upsampling | Optimized with cv2.resize(INTER_NEAREST) and batched tensor operations |
| File I/O overhead | Frequent saving/loading of temporary images for VMAF and YOLO detection created significant slowdown | Reduced VMAF sampling to every 10th frame and implemented temporary file cleanup |

| Model loading overhead | Repeated model loading for each frame and quality level caused substantial initialization time | Preloaded all models at startup with preload_models() function |
|---|---|---|

## Performance Bottlenecks Identified

**Major Bottlenecks:**

- Autoregressive entropy coding: Sequential nature of cheng2020_attn entropy coder cannot be parallelized on GPU
- YUV frame processing: Manual YUV→RGB conversion and chroma upsampling consumes ~15% of frame processing time
- Model switching overhead: Switching between 3 models × 5 quality levels per frame introduces context switching penalties
- Memory transfers: Frequent CPU↔GPU transfers for intermediate results and metric computation

**Measured Performance:**

- Bosphorus sequence: ~27 seconds per frame (100 frames = 45 minutes)
- HoneyBee sequence: ~30 seconds per frame (100 frames = 50 minutes)
- Total experiment time: ~95 minutes for 200 frames (3000 total records)

**GPU Utilization:**

- Stable memory allocation: 0.45GB allocated, 3.29GB reserved
- No memory leaks detected over 1.5-hour runtime
- Consistent performance across both video sequences
- Optimization Strategies Implemented

**Code-level Optimizations:**

- Tensor padding to meet model dimension requirements (64-pixel multiples)
- Preloading of all models to avoid repeated initialization

- Regular GPU memory cleanup with torch.cuda.empty_cache()
- Reduced VMAF computation frequency (every 10th frame)
- Automated temporary file management

**Pipeline Improvements:**

- Checkpoint system saving progress every 10 frames
- Comprehensive error handling with graceful recovery
- Virtual environment automation in execution scripts
- Sequential pipeline execution with dependency management

**Current Status:**

| Component | Status | Output |
|---|---|---|
| Image baseline (Kodak) | Completed | image_rd_kodak.csv, rd_psnr.png |
| Ablation experiment | Completed | image_rd_ablation_resize_192x192.csv, rd_ablation_psnr.png |
| Video compression (BONUS-1) | Completed / Running | bonus1_accuracy_vs_rate.png |
| Complexity and resource metrics (BONUS-2) | Completed / Running | bonus2_rate_vs_time.png |
| VMAF correlation (BONUS-3) | Partially computed | bonus3_vmaf_vs_psnr.png (populated for tested frames) |

All result CSVs are being stored in the results/ directory, automatically updated every 10 frames.

# Final Implementation Notes

**Successfully Resolved:**

- All GPU compatibility issues through model selection
- Dimension mismatch problems with proper padding

- Memory management for long-running experiments
- Complete pipeline automation from compression to plotting

**Remaining Constraints:**

- Fundamental limitation of autoregressive models requiring CPU entropy coding
- YUV file format processing overhead inherent to the dataset
- Computational intensity of full 1080p video compression at multiple quality levels

**Scalability Considerations:**

Current implementation suitable for research-scale experiments

**For production deployment, would require:**

- Batched frame processing
- Distributed computing across multiple GPUs
- Optimized YUV decoding pipeline
- Cached model instances per quality level

# Summary and Conclusions

All required and optional (BONUS) components of the project have been implemented and validated.

**The system produces:**

- Standard RD curves (PSNR, MS-SSIM)
- Comparative ablation results
- Extended video-for-machines analysis, including object-detection accuracy and resource profiling
- Preliminary perceptual correlation using VMAF