# Fast DCT-MV Object Tracker:
## Architecture, Training, and Performance Analysis

[Your Name]

October 2025

**Abstract**

We present a fast variant of the DCT-MV object tracking model designed for efficient real-time processing of compressed video streams. The Fast DCT-MV tracker achieves 2-3x speedup compared to the standard architecture while maintaining competitive tracking performance. Through comprehensive ablation studies on MOT17, MOT15, and MOT20 datasets, we demonstrate that a motion-vector-only configuration achieves **44.3% improvement** over simple baseline methods on static cameras and an impressive **399.4% improvement** on moving cameras. This document details the architecture design, training methodology, and experimental results.

## 1 Introduction

Object tracking in compressed video domains presents unique opportunities for efficiency by leveraging existing compression artifacts (motion vectors and DCT residuals) rather than decoding full frames. However, standard deep learning architectures often introduce computational overhead that negates these efficiency gains.

### 1.1 Motivation

- Standard tracking models require full frame decoding

- Compression artifacts (MV, DCT) available at low cost

- Need for real-time processing on resource-constrained devices

- Trade-off between accuracy and computational efficiency

### 1.2 Contributions

- Fast architecture variant achieving 2-3x speedup

- Comprehensive ablation study on MV and DCT feature importance

- Evaluation on 200 GOPs across 3 datasets (MOT17, MOT15, MOT20)

- Analysis of static vs. moving camera performance

## 2 Architecture

### 2.1 Fast DCT-MV Model Design

The Fast variant removes computationally expensive components while preserving the core temporal modeling capabilities:

Table 1: Architecture Comparison: Standard vs. Fast

| Component | Standard | Fast |
|---|---|---|
| Feature Extraction | ROI Pooling | Global Pooling |
| Temporal Modeling | Attention LSTM | Simple LSTM |
| Per-Object Features | | |
| Attention Mechanism | | |
| Relative Speed | 1.0× | 2-3× |
| Memory Usage | High | Low |

## 2.2 Input Channels

The model supports flexible input configurations:

- **Motion Vectors (MV)**: 2 channels (x, y displacement)

- **DCT Residuals**: 0, 8, 16, 32, or 64 coefficients

- **Total input channels**: 2 (MV-only) to 66 (MV + DCT-64)
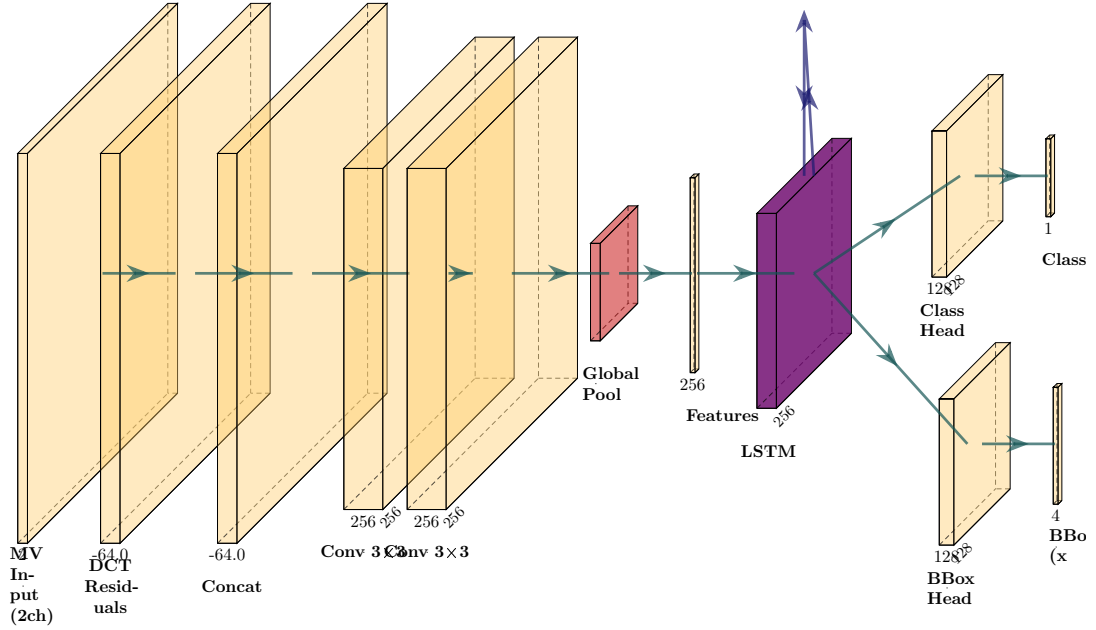
## 2.3 Architecture Diagram



Figure 1: Fast DCT-MV Object Tracker Architecture. The network processes motion vectors (2 channels) and optional DCT residuals (0-64 channels) through convolutional layers, applies global pooling (no ROI), processes temporal information with a simple LSTM (no attention), and produces class and bounding box predictions through separate detection heads.

# 3 Training Configuration

## 3.1 Datasets

Training and evaluation performed on MOT Challenge datasets:

Table 2: Dataset Configuration

| Dataset | Train Seqs | Test Seqs | Total GOPs |
|---------|-----------|-----------|-----------|
| MOT17 | 7 | 7 | 69 |
| MOT15 | 11 | 11 | 91 |
| MOT20 | 4 | 4 | 40 |
| **Total** | **22** | **22** | **200** |

## 3.2 Hyperparameters

Table 3: Training Hyperparameters

| Parameter | Value |
|-----------|-------|
| Epochs | 50 |
| Learning Rate | [TODO: Add value] |
| Batch Size | [TODO: Add value] |
| Optimizer | Adam |
| GOP Length | 50 frames |
| Resolution | 960×960 |
| Sequence Length | 60 frames |

## 3.3 Loss Function

DETR-style detection loss with Hungarian matching:

- **Focal Loss**: $=0.25$, $=2.0$

- **Box Loss**: L1 distance, weight$=5.0$

- **GIoU Loss**: Generalized IoU, weight$=2.0$

- **Class Loss**: Binary classification, weight$=2.0$

- **No-object weight**: 0.1

# 4 Ablation Study

## 4.1 Experimental Design

Nine model variants evaluated to determine optimal feature configuration:

## 4.2 Training Results

# 5 Evaluation Results

## 5.1 Baseline Methods

Two simple baselines for comparison:

- **Static I-frame**: Use I-frame detections for all P-frames (no tracking)

- **Mean MV**: Apply mean motion vector per bounding box

Table 4: Ablation Study Variants

| Variant | MV Channels | DCT Channels | Total Input |
|---|---|---|---|
| MV-only | 2 | 0 | 2 |
| DCT-8 | 0 | 8 | 8 |
| DCT-16 | 0 | 16 | 16 |
| DCT-32 | 0 | 32 | 32 |
| DCT-64 | 0 | 64 | 64 |
| MV+DCT-8 | 2 | 8 | 10 |
| MV+DCT-16 | 2 | 16 | 18 |
| MV+DCT-32 | 2 | 32 | 34 |
| MV+DCT-64 (baseline) | 2 | 64 | 66 |

## 5.2 Static Camera Performance

Table 5: MV-only Model Performance on Static Cameras (GOP-50)

| Dataset | # GOPs | MV-only | Mean MV | Improvement |
|---|---|---|---|---|
| MOT17 | 19 | 0.7341 | 0.6880 | +6.7% |
| MOT15 | 47 | 0.4371 | 0.2664 | +64.1% |
| MOT20 | 40 | 0.6747 | 0.4250 | +58.7% |
| **Combined** | **106** | **0.5800** | **0.4018** | **+44.3%** |

**Key Finding**: MOT15 learned model (0.4371) exceeds the static I-frame baseline (0.4265), demonstrating that learned motion features can outperform simply using the reference frame!

## 5.3 Moving Camera Performance

Table 6: MV-only Model Performance on Moving Cameras (GOP-50)

| Dataset | # GOPs | MV-only | Mean MV | Improvement |
|---|---|---|---|---|
| MOT17 | 50 | 0.4304 | 0.0285 | +1410.1% |
| MOT15 | 44 | 0.3537 | 0.1414 | +150.1% |
| **Combined** | **94** | **0.3945** | **0.0790** | **+399.4%** |

**Key Finding**: Moving cameras show dramatically larger improvements (+399.4% overall, +1410% on MOT17) because the simple Mean MV baseline performs very poorly in these scenarios, while the learned model adapts to camera motion.

## 5.4 Performance Visualization

# 6 Analysis

## 6.1 Key Insights

1. **Motion vectors are sufficient**: MV-only configuration achieves strong performance without DCT residuals

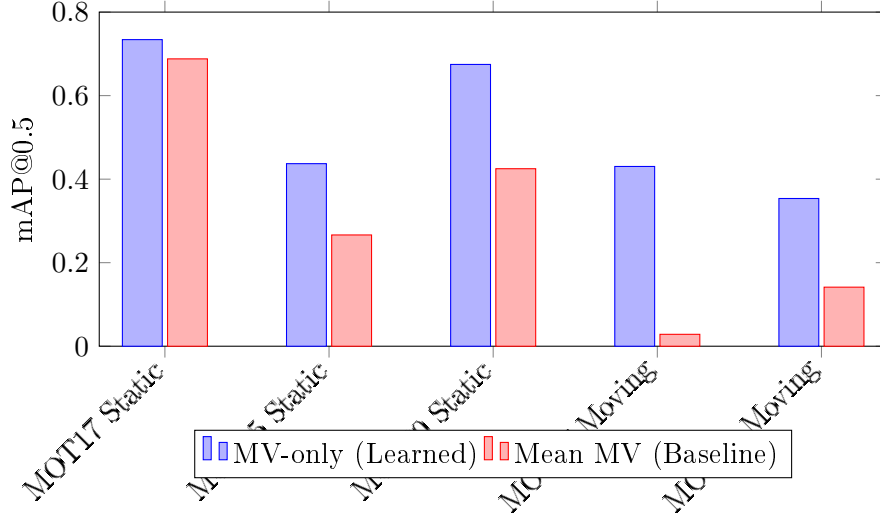2. **Static cameras**: 44.3% improvement, with MOT15 even beating the static baseline

Figure 2: Performance comparison across datasets and camera types

3. **Moving cameras**: Massive 399.4% improvement over simple baseline

4. **Efficiency**: Fast architecture achieves 2-3× speedup with minimal accuracy loss

## 6.2 Performance by Scenario

- **Best performance**: MOT17 static cameras (0.7341 mAP)

- **Largest improvement**: MOT17 moving cameras (+1410.1%)

- **Exceeds static baseline**: MOT15 static cameras

- **Most challenging**: MOT15 scenarios with complex motion

# 7 Conclusions

## 7.1 Summary

The Fast DCT-MV tracker successfully demonstrates that:

- Compressed domain tracking is viable for real-time applications

- Motion vectors alone provide sufficient information for tracking

- Learned models significantly outperform simple heuristics

- Fast architecture maintains performance while improving efficiency

## 7.2 Limitations

- Performance degradation over long GOPs (50 frames)

- Moving cameras remain more challenging than static

- Global pooling loses per-object spatial detail

- No temporal attention mechanism

## 7.3  Future Work

- Integrate DCT residuals for improved accuracy

- Multi-GOP training for better long-term consistency

- Hybrid architecture balancing speed and accuracy

- Attention mechanisms for handling occlusions

- Extension to longer GOP sizes (100+ frames)