

FluxShard: Distributed Motion-Aware Cache Remapping for Real-Time Video Analytics

IEEE Publication Technology Department

Abstract—Edge-cloud video analytics caches intermediate features across consecutive frames to cut redundant computation and transmission. Existing methods operate at whole-scene granularity—reusing or invalidating entire feature maps—and assume a roughly stationary scene. Even moderate motion misaligns the cache with the current frame, triggering widespread invalidation although most content has merely shifted rather than truly changed. The root cause is a *granularity mismatch*: the cache is treated as an indivisible unit, yet real-world motion is local and heterogeneous.

We present FluxShard, a motion-aware edge-cloud analytics system that elevates codec-level motion vectors (MVs) into a first-class control signal for feature caching. FluxShard decomposes the cached feature map into block-level shards that flow with observed motion, realigning reusable content and isolating genuinely changed regions as a minimal recomputation set. A *Receptive-Field Alignment Principle* guarantees bit-equivalent correctness of MV-guided reuse across convolutional layers, and a *profiling-driven truncation policy* sustains high sparsity under growing motion with negligible accuracy loss. At runtime, FluxShard adaptively routes the sparse residual workload between edge and cloud based on network conditions and device load. Evaluation on real-world dynamic video sequences spanning object detection and instance segmentation shows that FluxShard achieves up to 92% bandwidth reduction, 5.5× speedup, and ~99% accuracy retention over state-of-the-art baselines.

Index Terms—Video analytics, robotics, edge-cloud collaborative system, CNN inference acceleration

I. INTRODUCTION

Video analytics underpins latency-critical applications such as embodied intelligence [25], [10], aerial drones [1], [19], augmented reality [26], and smart surveillance [13]. Processing solely on edge devices avoids transmission delay but is constrained by limited compute and energy; offloading to the cloud provides ample resources but introduces network latency. A natural middle ground is to *reuse* previously computed features across consecutive frames: since adjacent frames share substantial visual content, caching and reusing unchanged results reduces both redundant computation and transmission.

Existing reuse methods, however, operate at the granularity of the *whole scene* and are therefore brittle under motion. Pipeline-level approaches [16], [8] cache entire intermediate feature maps and reuse them only when successive frames are globally similar. Delta-based methods [2], [18], [22] maintain a scene-level reference and update it via pixel-wise differences; MotionDeltaCNN [21] extends this with a single global homography. Because all these approaches treat the cache as an indivisible entity, any camera ego-motion or object movement triggers widespread cache invalidation—even when

most content has merely *shifted* rather than changed. The root cause is a *granularity mismatch*: the cache is monolithic, yet real-world motion is local and heterogeneous.

Block-level motion vectors (MVs) from standard video codecs offer a natural remedy. Extracted as a free byproduct of decoding, MVs capture per-region displacement between consecutive frames and can therefore distinguish spatial shift from genuine content change. Re-indexing cached features according to these displacements recovers reusable content that whole-scene methods would discard, while regions where MVs are inherently insufficient—non-rigid deformation, disocclusion, newly appearing content—are isolated as a minimal *residual set* that truly requires recomputation.

Realizing this idea, however, introduces two challenges. (1) the **receptive field inconsistency** problem: even when cached features are correctly aligned via MVs, directly reusing the aligned output of a layer that aggregates over a spatial neighborhood of input (i.e., receptive field) with size larger than one is not always valid, because heterogeneous motion can assemble a neighborhood that differs from the one used to produce the cached output, silently corrupting results. Naïvely detecting such mismatches requires checking the full receptive field at every output position of every layer, incurring prohibitive cost rivaling recomputation itself. (2) the **sparsity decay** problem: the residual set tends to expand through cascaded layers as layers with receptive field size larger than one causes recomputed positions to bleed into their neighbors, progressively eroding the efficiency advantage of selective recomputation as motion intensity grows.

We propose **FluxShard**, a motion-aware edge-cloud video analytics system that treats the feature cache as block-level *shards* flowing freely with motion vectors. FluxShard uses MVs as a first-class control signal to align cached features with the current frame, isolate the minimal residual set, and dynamically route work—local sparse update on the edge or selective offloading of only residual regions to the cloud—based on residual size, network conditions, and device load.

To address receptive field inconsistency, FluxShard introduces the **Receptive Field Alignment Principle (RFAP)** (§IV-B), which folds the per-layer, per-position neighborhood check into a single lightweight pass over the input-level MV field. RFAP produces a conservative reuse mask at each layer in time linear in the input resolution, guaranteeing bit-equivalent correctness without per-layer verification overhead.

To counter sparsity decay, FluxShard employs a **profiling-driven truncation policy** (§IV-C). Offline, we characterize the tightest accuracy-preserving truncation thresholds across a range of motion intensities; at runtime, the system indexes

Grouped bars: edge-only latency vs. cloud RTT
at uplink = 5, 10, 20, 50 Mbps

Figure 1. Per-frame latency under varying uplink bandwidth. Neither pure-edge nor pure-cloud meets real-time across all conditions.

Line plot: reuse ratio (%) vs. avg motion (px/frame)
for a whole-scene cache baseline

Figure 2. Cache reuse ratio versus per-frame motion. Whole-scene granularity causes reuse to collapse under moderate motion.

into this mapping based on observed motion, dynamically modulating truncation to sustain high sparsity with negligible accuracy loss.

In summary, this paper makes the following contributions:

- We identify granularity mismatch as the key limitation of whole-scene cache reuse and show that codec-level motion vectors provide the right abstraction to overcome it.
- We propose FluxShard, a system that decomposes the feature cache into motion-aligned shards and co-optimizes sparse recomputation with adaptive edge-cloud workload routing.
- We establish the Receptive Field Alignment Principle for bit-equivalent MV-guided feature reuse across layers, and a profiling-driven truncation policy that sustains high sparsity under varying motion with bounded accuracy loss.
- We evaluate FluxShard on object detection and instance segmentation with YOLOv11 [15] on two dynamic video benchmarks (DAVIS [23], 3DPW [24]). Compared to the strongest baseline, FluxShard achieves up to $5.5\times$ end-to-end speedup and 92% bandwidth reduction while retaining over 99% of full-model accuracy.

II. CHALLENGES AND MOTIVATION

A. The Edge-Cloud Dilemma

Edge-cloud video analytics must navigate two competing bottlenecks. Running YOLOv11m [15] on an NVIDIA Jetson Xavier NX takes **XX** ms per 1080p frame, far exceeding a 33 ms real-time budget. Offloading every frame removes the compute bottleneck but introduces a transmission one: JPEG-compressed 1080p at 30 fps sustains **XX** Mbps, routinely exceeding edge uplink capacity. Figure 1 quantifies this tension.

Feature cache reuse bridges this gap: when consecutive frames share content, cached results substitute for fresh computation. However, existing mechanisms [16], [8], [22], [21] treat the cache as an indivisible whole-scene entity—whether the decision is binary (reuse or recompute the entire frame) or pixel-level (propagate a dense difference map). This granularity breaks down under motion. Figure 2 shows that even modest displacement triggers widespread invalidation: the reuse ratio drops from **XX**% on near-static sequences to below **XX**% once motion exceeds **XX** px/frame, despite most content being merely shifted rather than truly changed.

Bars: residual set size w/o vs. w/ MV alignment
across low / medium / high motion bins

Figure 3. MV alignment eliminates displacement-induced false misses, cutting the residual set by **XX-XX**%.

Bars: mAP for oracle / whole-scene reuse / naïve MV
reuse

Figure 4. Naïve MV-aligned reuse destroys accuracy due to receptive field inconsistency, even when the residual set is small.

Takeaway. Cache reuse is essential, but whole-scene granularity cannot tolerate motion. A finer-grained, motion-aware mechanism is needed.

B. Motion Vectors: Opportunity and Correctness Failure

H.264/H.265 codecs estimate a per-block displacement (d_x, d_y) as part of normal encoding, providing per-region motion information at *zero additional cost*. While recent studies [7], [6], [4] have shown that the codec’s full reconstruction pipeline—transform, quantization, etc—distorts feature statistics in ways that degrade downstream inference accuracy, the motion vectors themselves are a pure geometric signal untouched by these lossy stages. By shifting cached feature blocks according to these MVs, displaced content is realigned before differencing, and only genuinely changed regions—disocclusion, deformation, new objects—remain in the *residual set*. Figure 3 confirms MV alignment reduces the residual set by **XX-XX**% across motion intensities: most whole-scene cache misses stem from displacement, not true change.

Despite the smaller residual set, directly reusing MV-aligned features produces catastrophic accuracy loss. Figure 4 shows mAP drops from **XX**% to **XX**% under naïve MV reuse. The root cause is *receptive field inconsistency*: layers with receptive field size > 1 aggregate spatial neighborhoods. When adjacent blocks carry different MVs, re-indexing assembles patches that were never contiguous in the original frame; the cached output was computed from a different neighborhood. The resulting error compounds through cascaded layers.

Takeaway. MV alignment eliminates most false cache misses, but exploiting it requires a correctness guarantee under heterogeneous per-block motion. This is the *first challenge* our design must address.

C. Sparsity Decays Through Cascaded Layers

Even with correct reuse, a second challenge erodes efficiency. Each convolutional layer with kernel radius r reads r positions beyond every active position, expanding the residual set by a margin of r per side per layer. Over L layers the expansion compounds: a single active input position inflates to a region of radius $\sum_l r_l$ at the final layer.

Figure 5 traces per-layer residual ratio through YOLOv11m. A modest input residual ratio of **XX**% grows past **XX**% by mid-network and saturates near 100% before the last layer.

Lines: per-layer residual ratio (%) vs. layer index
curves for input residual ratio = 5, 10, 20, 30%

Figure 5. Receptive field bleed causes the residual set to expand at each layer, approaching full recomputation before the final layer.

At that point selective recomputation offers no advantage over full inference.

Takeaway. Without an explicit mechanism to arrest receptive field bleed, selective recomputation degenerates into full recomputation under real-world motion. A sparsity-preserving strategy with bounded accuracy loss is the *second challenge* our design must address.

III. SYSTEM OVERVIEW

This section presents the FluxShard system. We first define the system model and optimization objective (§III-A), then describe the end-to-end pipeline that addresses the two challenges identified above (§III-B).

A. System Model

We consider an edge-cloud video analytics system in which an edge device (e.g., an NVIDIA Jetson) receives a live video stream from a co-located camera and a remote cloud server provides supplementary compute capacity. The two endpoints communicate over a bandwidth-limited uplink with round-trip latency that varies over time. Each endpoint hosts an identical copy of the inference model and maintains a per-layer feature cache that stores the output produced during its most recent inference. For each incoming frame, the system makes a *dispatch decision*: whether to run inference locally at the edge or offload it to the cloud.

a) Notation.: Let $\{I_t\}$ denote the incoming frame sequence. We write \mathbf{m}_t for the block-level motion vector (MV) field extracted from the codec when decoding frame I_t ; each entry maps a block position in I_t to its corresponding position in the reference frame.

The DNN \mathcal{N} comprises L layers indexed by $l \in \{1, \dots, L\}$. Layer l takes an input feature map $\mathbf{F}^l \in \mathbb{R}^{H_l' \times W_l' \times C_l'}$ and produces an output feature map $\mathbf{O}^l \in \mathbb{R}^{H_l \times W_l \times C_l}$; by convention $\mathbf{F}^1 = I_t$. Since our analysis concerns only spatial positions, we omit the channel dimension hereafter. The output of a layer at spatial position (i, j) depends on a set of input positions called its *receptive field*, denoted $\mathcal{R}^l(i, j) \subseteq \{1, \dots, H_l'\} \times \{1, \dots, W_l'\}$, with radius r_l . For convolutional and linear layers, \mathbf{w}^l denotes the weight tensor applied over this receptive field.

Each endpoint caches the output of every layer from its most recent inference; we write $\hat{\mathbf{O}}^l$ for the cached output at layer l and $\hat{\mathbf{F}}^l$ for the corresponding cached input. Given the MV field, position (i, j) in the current frame maps to a cached position (\hat{i}, \hat{j}) ; likewise, each input position $(p, q) \in \mathcal{R}^l(i, j)$ maps to its cached counterpart (\hat{p}, \hat{q}) .

A *shard* is a block-sized unit of a feature map whose spatial extent corresponds to one entry in the MV field. At deeper layers with reduced spatial resolution, neighboring MVs that

map to the same shard are aggregated and their magnitudes divided by the stride factor, so that each shard retains a one-to-one correspondence with a single displacement measured in feature-map coordinates.

b) Assumptions.: The inference result (e.g., detections, segmentation masks) is consumed at the executing endpoint or forwarded to a co-located downstream service. Downstream processing is typically lightweight (e.g., issuing an alert, logging an action label), and the processed outcome is negligible in size; we therefore treat result delivery as outside the latency-critical path and exclude it from our optimization scope.

c) Objective.: Let $d_t \in \{\text{edge}, \text{cloud}\}$ denote the dispatch decision for frame I_t . The per-frame end-to-end latency $T(d_t)$ comprises local computation, data transmission, and remote computation as applicable. Let $\mathcal{A}(d_t)$ denote the task accuracy achieved under decision d_t and \mathcal{A}^* the accuracy of full recomputation without caching. In general, latency reduction techniques trade off accuracy for speed, so $\mathcal{A}(d_t)$ depends on the specific optimizations applied at the selected endpoint. Our objective is:

$$\min_{d_t} T(d_t) \quad \text{s.t.} \quad \mathcal{A}(d_t) \geq \alpha \mathcal{A}^*. \quad (1)$$

B. Pipeline

Figure 6 illustrates the end-to-end FluxShard pipeline. Each endpoint maintains a *dispatch layer*: a lightweight book-keeping structure that stores (i) an *input cache*—the input used in the endpoint’s most recent inference—and (ii) an *accumulated MV field* from that cached input to the present frame. The dispatch layer enables per-endpoint reusability estimation without running the model.

When an encoded frame arrives at the edge, FluxShard proceeds in four stages.

a) Stage 1: MV extraction.: Block-level motion vectors \mathbf{m}_t are extracted from the video codec as a free byproduct of decoding, capturing per-region motion at no additional cost. Both dispatch layers incorporate \mathbf{m}_t into their respective accumulated MV fields.

b) Stage 2: Reusability estimation.: Each dispatch layer uses its accumulated MV field to shift its cached input, aligning displaced content with the current frame. A per-block comparison between the aligned cache and the decoded frame identifies positions with non-negligible difference; these positions form the *residual set*. From the residual set size, the edge dispatch layer estimates local sparse inference latency, while the cloud dispatch layer estimates the cost of transmitting the MV field and residual pixels plus remote inference under current network conditions.

c) Stage 3: Dispatch decision.: The frame is assigned to the endpoint with the lower estimated latency, implicitly adapting to motion complexity, network bandwidth, and endpoint load.

d) Stage 4: Inference and cache update.: The selected endpoint applies its accumulated MV field to rearrange cached shards across all layers and fills in residual pixels at the input, producing the updated input; its input cache is then replaced and its accumulated MV field reset. The residual set is propagated layer by layer through the model: at each layer,

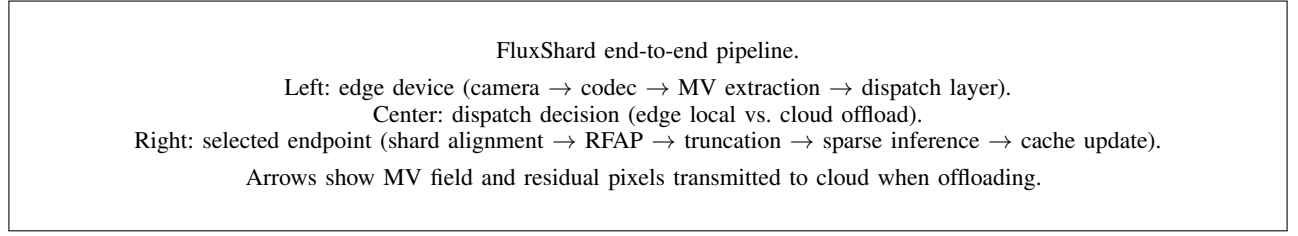


Figure 6. Overview of the FluxShard pipeline. Numbered stages correspond to the description in §III-B.

the Receptive Field Alignment Principle (§IV-B) identifies positions that can be safely reused despite heterogeneous motion, and a truncation policy (§IV-C) suppresses positions whose residual magnitude is sufficiently small. Only the remaining positions are recomputed; their outputs are merged into the layer’s feature cache. The final result is produced by fusing reused and freshly computed features at the last layer.

The non-selected endpoint retains its accumulated MV field, allowing its cache to span a longer interval for future reuse. When the server is selected, the edge transmits only the accumulated MV field and the residual pixels; a lightweight *receive layer* on the server mirrors the dispatch layer logic to reconstruct the full input before inference.

IV. SYSTEM DESIGN

The pipeline in §III-B described the four stages that every frame traverses; this section develops the mechanisms deferred in that overview. We first formalize when a cached output, after MV-guided shard alignment, can be directly reused without accuracy loss (§IV-A); this criterion underpins all transmission and computation reductions of FluxShard while preserving inference accuracy. We then derive the Receptive Field Alignment Principle (§IV-B), which makes this criterion efficiently evaluable under per-block heterogeneous motion by folding multi-layer consistency checks into a single pass over the input-level MV field.

Next, we introduce a profiling-driven truncation policy (§IV-C) that compensates for the sparsity decay caused by receptive field expansion across cascaded layers, maintaining high computational savings with bounded accuracy loss. Finally, we present the dispatch scheduler (§IV-D) that integrates the reuse map, current network conditions, and endpoint load into a per-frame decision of whether to update residual regions locally at the edge or offload them to the cloud, minimizing end-to-end latency while satisfying the accuracy constraint in Eq. (1).

A. Layer-Local Reuse Criterion

Stage 4 of the pipeline (§III-B) propagates the residual set layer by layer, reusing cached outputs wherever possible. The quality of the final result hinges on a per-position decision: after MV-guided shard alignment, is the cached output close enough to the true output to be reused? This subsection formalizes that decision as a layer-local criterion and identifies the structural obstacle that heterogeneous MVs pose for its efficient evaluation.

a) *Setup and notation.*: Three quantities coexist at every output position (i, j) of layer l :

- $\mathbf{O}^l(i, j)$: the *ground-truth* output, obtained by applying layer l to the current input \mathbf{F}^l ;
- $\hat{\mathbf{O}}^l(i, j)$: the *cached* output from a previous frame, where $(\hat{i}, \hat{j}) = (i, j) - \mathbf{d}_{\text{out}}^l(i, j)$ is the MV-aligned source position;
- $\tilde{\mathbf{O}}^l(i, j)$: the *assembled* output actually consumed by layer $l+1$.

b) *Reuse validity.*: Reuse at position (i, j) is *valid* when the substitution of the cached value for the ground-truth introduces bounded error:

$$|\mathbf{O}^l(i, j) - \hat{\mathbf{O}}^l(i, j)| \leq \tau, \quad (2)$$

where τ is a task-driven tolerance; $\tau = 0$ recovers exact reuse. This condition is universal across layer types and serves as the invariant maintained across the full layer stack.

c) *Layer output assembly.*: Based on this criterion, FluxShard assembles the output of layer l as

$$\tilde{\mathbf{O}}^l(i, j) = \begin{cases} \hat{\mathbf{O}}^l(i, j), & \text{if Eq. (2) holds,} \\ \mathbf{O}^l(i, j), & \text{otherwise.} \end{cases} \quad (3)$$

Eq. (3) formalizes the core operation of FluxShard: the assembled output is a mosaic of *MV-aligned cached values* at reusable positions and *freshly computed values* at residual positions, and serves as the input to layer $l+1$. Maximizing the reusable set while respecting Eq. (2) is the objective of the remainder of this section.

d) *Input-side sufficient condition.*: Eq. (2) cannot be checked directly without computing $\mathbf{O}^l(i, j)$, which defeats the purpose. For a linear layer, linearity yields

$$|\mathbf{O}^l(i, j) - \hat{\mathbf{O}}^l(i, j)| \leq \|\mathbf{w}^l\|_1 \cdot \max_{(p, q) \in \mathcal{R}^l(i, j)} |\mathbf{F}^l(p, q) - \hat{\mathbf{F}}^l(\hat{p}, \hat{q})|, \quad (4)$$

so Eq. (2) is guaranteed whenever the input patch satisfies

$$\max_{(p, q) \in \mathcal{R}^l(i, j)} |\mathbf{F}^l(p, q) - \hat{\mathbf{F}}^l(\hat{p}, \hat{q})| \leq \frac{\tau}{\|\mathbf{w}^l\|_1}. \quad (5)$$

For nonlinear activations that are 1-Lipschitz (ReLU, sigmoid), the input-side bound directly implies the output-side one, so no extra margin is needed. Crucially, Eq. (5) is evaluated using only layer l ’s own input and cached output; it never recurses to the model input.

This avoids computing \mathbf{O}^l itself, yet the check still inspects every position in the receptive field $\mathcal{R}^l(i, j)$ at layer l ’s input, whose cost is comparable to executing layer l . To make the scheme practical, we next show how to propagate the reuse

decision from a single pass over the *model input*, eliminating the need for per-layer inspection entirely.

e) Free propagation through pointwise layers.: When $|\mathcal{R}^l| = 1$ (1×1 convolution, pointwise activation), the reuse decision propagates without re-evaluation: reused positions carry zero input difference and satisfy any threshold, while recomputed positions carry nonzero difference that cannot pass the tighter bound $\tau / \|\mathbf{w}^{l+1}\|_1 \leq \tau$. The reuse map is therefore inherited unchanged by every subsequent pointwise layer at the same spatial resolution.

f) Breakdown under spatial aggregation.: Once $|\mathcal{R}^l| > 1$ (e.g., 3×3 convolution), free propagation breaks down. Each output aggregates a spatial neighborhood whose elements may have been relocated from *disjoint* cached positions by heterogeneous motion vectors. Even if every individual input position was reused with zero difference, the spatial composition within the receptive field can differ from the one that produced the cached output, invalidating Eq. (5) and forcing a per-position re-evaluation of the criterion at every such layer. As shown in Fig. 4, ignoring this inconsistency causes accuracy to degrade sharply with motion magnitude, dropping over XX% under moderate motion.

Delta methods [22], [9], [21] avoid this expense by restricting all blocks to a single shared MV, under which the reuse map can be seeded at the model input and propagated by receptive field inclusion alone; however, as shown in Fig. 3, this uniform-MV assumption collapses under heterogeneous motion. No prior work addresses the combination of per-block heterogeneous MVs with layer-wise feature cache reuse, where the consistency problem is inherent. The next subsection derives conditions that restore efficient evaluation without restricting the MV field.

B. Receptive Field Alignment Principle

The preceding subsection showed that spatial aggregation layers invalidate free propagation of the reuse map: each output position aggregates a neighborhood whose elements may originate from disjoint cached regions under heterogeneous motion. Checking the reuse criterion at every such layer costs as much as recomputation. This subsection derives the *Receptive Field Alignment Principle* (RFAP), which identifies two structural conditions on the input-level MV field that, when jointly satisfied, guarantee exact reuse at all downstream layers without per-layer inspection.

1) Two Conditions for Correct Reuse: The above conflict arises because each layer l carries two independent MV fields— \mathbf{d}_{in}^l on its input grid and $\mathbf{d}_{\text{out}}^l$ on its output grid—and heterogeneous motion can make them geometrically incompatible. Receptive Field Alignment Principle (RFAP) imposes two conditions that are jointly sufficient for reuse at any output position (i, j) ; violating either one can produce silent errors that accumulate across layers.

Condition 1: Intra-receptive-field uniformity. All input positions within the receptive field must share a single displacement:

$$\mathbf{d}_{\text{in}}^l(p, q) = \mathbf{d}_{\text{in}}^l(p', q'), \quad \forall (p, q), (p', q') \in \mathcal{R}^l(i, j). \quad (6)$$

A uniform shift preserves the relative arrangement of values consumed by the layer. If violated, the aligned patch stitches together fragments from non-contiguous cached regions—an input composition that never existed in the cached frame.

Condition 2: Input-output geometric coherence. Let $\pi_{p,q}^l(i, j) = s^l \cdot (i, j) + (p, q)$ map output position (i, j) to the input coordinate of receptive field entry (p, q) , where s^l is the stride of layer l . When projected into input coordinates, the output-level displacement must land on the same cached position as the input-level displacement at every receptive field entry:

$$\pi_{p,q}^l((i, j) + \mathbf{d}_{\text{out}}^l(i, j)) = \pi_{p,q}^l(i, j) + \mathbf{d}_{\text{in}}^l(\pi_{p,q}^l(i, j)), \quad \forall (p, q) \in \mathcal{R}^l(i, j). \quad (7)$$

For the affine $\pi_{p,q}^l$ above, this simplifies to $\mathbf{d}_{\text{out}}^l = \mathbf{d}_{\text{in}}^l / s^l$, linking the two fields through the layer stride. If violated, the cached output and the aligned input patch refer to inconsistent cached content: the two alignment paths diverge.

When both conditions hold, every input to the layer is identical between the current and cached evaluations. By determinism of the layer, $\mathbf{O}^l(i, j) = \hat{\mathbf{O}}^l(\hat{i}, \hat{j})$, satisfying Eq. (2) with $\tau = 0$, i.e., exact reuse. Conversely, violating either condition can cause silent corruption that accumulates across layers. The set of output positions that fail either condition at layer l constitutes the *MV inconsistency mask* \mathcal{M}^l ; these positions must be recomputed.

a) Practical scope.: Enforcing both conditions at every spatial-aggregation layer would negate the sparsity gains that motivate reuse. The two conditions, however, differ sharply in their impact on sparsity. Condition 1 must be checked at every layer whose receptive field exceeds one; such layers typically constitute a large share of the backbone (e.g., XX% to YY% in ResNet-50, EfficientDet and YOLOv11). Under strong motion, neighboring positions frequently carry inconsistent MVs, so strict enforcement at every such layer marks a large fraction of positions as inconsistent and collapses the residual set toward full recomputation. Condition 2, by contrast, is relevant only at layers that change spatial resolution; these often number just several in a typical backbone (e.g., 3–4 in ResNet-50, EfficientDet and YOLOv11), so its impact on overall sparsity remains limited.

We therefore enforce Condition 1 only at the first layer whose receptive field exceeds one, recomputing all inconsistent positions so that their outputs reflect the true assembled input. Condition 2 is enforced at every resolution-changing layer. Beyond the first spatial-aggregation layer, unchecked Condition 1 violations may weaken the per-layer bound in Eq. (2), but the relationship between per-layer feature error and end-task accuracy is inherently indirect: the bound serves as an operational proxy, not a formal accuracy guarantee. The actual accuracy guarantee comes from the profiling-driven truncation policy (§IV-C), which selects τ by directly measuring end-task accuracy under all approximations—including relaxed Condition 1 enforcement—and ensures $A(\tau) \geq \alpha \cdot A^*$. We verify in §?? that this strategy keeps accuracy loss below XX% while recovering XX% of the sparsity lost to strict enforcement.

C. Profiling-Driven Truncation Policy

After the first spatial-aggregation layer, positions in the residual set arise from two distinct sources: those violating RFAP and those whose underlying content has genuinely changed. Both sources produce a spectrum of discrepancy magnitudes: high-magnitude entries that matter for accuracy and low-magnitude entries that can be safely discarded.

With $\tau = 0$ every nonzero left-hand side of Eq. 5 triggers recomputation, preserving bit-exact correctness but collapsing sparsity toward the full-recomputation regime observed in §II-C. Raising τ allows positions whose input-patch discrepancy is small to be retained from cache, trading a controlled approximation for substantially higher sparsity.

Because the check in Eq. (5) operates on discrepancy magnitude regardless of origin, a single τ filters both sources uniformly: positions with large discrepancy are recomputed while those with small discrepancy are reused. The per-layer bound in Eq. (2) may not hold strictly at layers where Condition 1 is relaxed, but the relationship between per-layer feature error and end-task accuracy is inherently indirect. Rather than relying on this bound, the truncation policy selects τ by directly measuring end-task accuracy under the full pipeline—including RFAP relaxation, truncation, and cache assembly—so that all sources of approximation are absorbed into a single, empirically validated knob.

Formulation. The truncation operator retains the cached value at position p in layer ℓ whenever the left-hand side of Eq. 5 falls below $\tau/\|\mathbf{w}^\ell\|_1$, removing p from the residual set \mathcal{R}_ℓ . Because Eq. 5 normalizes the raw input-patch discrepancy by the layer-specific factor $\|\mathbf{w}^\ell\|_1$, the threshold τ is a single global scalar shared across all layers. The optimization objective is:

$$\max_{\tau \geq 0} \tau \quad \text{s.t.} \quad A(\tau) \geq \alpha \cdot A^*, \quad (8)$$

where $A(\tau)$ is end-task accuracy under threshold τ and A^* is the accuracy of full inference without any reuse.

Offline profiling. We solve (8) offline on a representative calibration set. For a discrete set of motion intensity levels—quantified by a scalar statistic extracted from the MV field—we sweep τ upward from zero and record the accuracy at each value. The output is a lookup table mapping each motion intensity level to the maximum feasible τ satisfying the accuracy constraint.

Empirically, the feasible τ increases with motion intensity: when motion is large, discrepancy magnitudes are collectively elevated, so a higher threshold still retains sufficient signal. This trend is favorable because aggressive truncation becomes available precisely in the high-motion regime where the residual set is largest and sparsity recovery is most needed.

Runtime selection. At each frame, the dispatch layer computes the motion statistic from the current MV field and looks up τ from the profiled table. At each layer ℓ , the left-hand side of Eq. 5 is evaluated per position: a position whose value falls below $\tau/\|\mathbf{w}^\ell\|_1$ is removed from \mathcal{R}_ℓ . The cost is one scalar lookup plus one element-wise comparison per layer, adding negligible overhead.

Following DeltaCNN [22], we restrict truncation to activation layers only. Because activations are the natural sparsification points in a DNN pipeline, applying the threshold there is sufficient to shrink the residual set; all other layers simply operate on whichever residual set they inherit. Furthermore, the threshold check needs to be evaluated only over positions that are currently *in* the residual set: positions already marked as reusable by RFAP or by alignment have zero discrepancy by construction and satisfy any $\tau > 0$ automatically.

Dispatch layer integration. The dispatch layer is an identity mapping prepended to the model, so Eq. (5) applies with $\|\mathbf{w}\|_1 = 1$: each MV-aligned input position is compared against its cached counterpart, and positions whose discrepancy falls below τ are excluded from the residual set. When inference is offloaded to the cloud, only the residual positions and the MV field are transmitted; the cloud reconstructs reusable positions from its own cache via the same alignment and merges them with the received residual positions to assemble the complete input for subsequent layers.

D. Dispatch Scheduler

The preceding subsections determine *what* to reuse and *what* to recompute; this subsection decides *where* the recomputation takes place. For each frame, the scheduler makes a binary decision—edge local inference or cloud offload—to minimize end-to-end latency while the accuracy constraint is already satisfied by the truncation policy.

a) *Cost estimation:* The edge maintains two dispatch layers, one tracking the local inference state and one mirroring the cloud inference state (§III). For each incoming frame, both dispatch layers independently perform shard alignment and residual set computation using their respective cached inputs and accumulated MV fields. Each dispatch layer derives its own motion statistic m_e, m_c and looks up the corresponding truncation threshold τ_e, τ_c from the profiled table, yielding residual sets \mathcal{S}_0^e and \mathcal{S}_0^c .

Estimating the end-to-end sparse inference latency from $|\mathcal{S}_0|$ alone is unreliable: the input-layer residual count does not capture how sparsity evolves across layers, which depends on the magnitude distribution of residual activations and the truncation threshold. We observe empirically that motion intensity subsumes these factors—higher motion produces both larger and higher-magnitude residual sets—making it a strong proxy for end-to-end latency. Accordingly, we profile sparse inference latency offline on each endpoint as a function of motion intensity and store the result as a lookup table (key: discretized motion intensity; value: measured latency). At runtime, edge latency is estimated as $T_{\text{edge}} \leftarrow \mathcal{T}_{\text{edge}}[m_e]$. Cloud latency combines the profiled cloud-side inference cost with a transmission term: $T_{\text{cloud}} \leftarrow \mathcal{T}_{\text{cloud}}[m_c] + |\mathcal{S}_0^c|/\hat{B}$, where \hat{B} is an exponentially weighted moving average of recent bandwidth measurements. The scheduler selects the path with lower total latency; when the two estimates are within a margin ϵ , it prefers cloud offload to conserve edge energy.

Algorithm 1 summarizes the complete per-frame pipeline, integrating MV extraction, reusability estimation, dispatch, and cache maintenance.

Algorithm 1: FluxShard Per-Frame Pipeline

Input: Current frame I_t ; edge dispatch layer (\hat{I}_e, \mathbf{d}_e);
cloud dispatch layer (\hat{I}_c, \mathbf{d}_c); profiled table \mathcal{T}_τ ;
latency tables $\mathcal{T}_{\text{edge}}, \mathcal{T}_{\text{cloud}}$

Output: Inference result \mathbf{y}_t

Extract MV field \mathbf{d}_t from codec for frame I_t ;
 $\mathbf{d}_e \leftarrow \text{Accumulate}(\mathbf{d}_e, \mathbf{d}_t)$;
 $\mathbf{d}_c \leftarrow \text{Accumulate}(\mathbf{d}_c, \mathbf{d}_t)$;
 $m_e \leftarrow \text{MotionStat}(\mathbf{d}_e)$;
 $m_c \leftarrow \text{MotionStat}(\mathbf{d}_c)$;
 $\tau_e \leftarrow \mathcal{T}_\tau[m_e]$; $\tau_c \leftarrow \mathcal{T}_\tau[m_c]$;
 $\mathcal{S}_0^e \leftarrow \text{ResidualSet}(I_t, \hat{I}_e, \mathbf{d}_e, \tau_e)$;
 $\mathcal{S}_0^c \leftarrow \text{ResidualSet}(I_t, \hat{I}_c, \mathbf{d}_c, \tau_c)$;
 $T_{\text{edge}} \leftarrow \mathcal{T}_{\text{edge}}[m_e]$;
 $T_{\text{cloud}} \leftarrow \mathcal{T}_{\text{cloud}}[m_c] + |\mathcal{S}_0^c| / \hat{B}$;
if $T_{\text{edge}} < T_{\text{cloud}} - \epsilon$ **then**
 $\mathbf{y}_t \leftarrow \text{SparseInfer}_{\text{edge}}(I_t, \mathcal{S}_0^e, \tau_e, \mathbf{d}_e)$;
 Update \hat{I}_e ; reset \mathbf{d}_e ;
else
 Transmit \mathcal{S}_0^c and \mathbf{d}_c to cloud;
 $\mathbf{y}_t \leftarrow \text{SparseInfer}_{\text{cloud}}(\mathcal{S}_0^c, \tau_c, \mathbf{d}_c)$;
 Update \hat{I}_c ; reset \mathbf{d}_c ;
return \mathbf{y}_t ;

V. RELATED WORK

Because consecutive video frames share substantial visual content, a growing body of work caches previously computed features and reuses them for subsequent frames, reducing both computation and transmission. Existing approaches fall into two broad categories.

Pipeline-level caching. Methods such as SPINN [16] and COACH [8] cache intermediate feature maps or label-level predictions and reuse them when successive inputs are deemed globally similar. COACH, for example, maintains semantic cluster centers in the feature space and triggers an early exit—bypassing cloud inference entirely when a new frame’s embedding falls within a similarity threshold. These methods make a *binary, whole-input* reuse decision: either the entire cached result is reused or it is fully recomputed. This coarse granularity suits image classification, where a single label summarizes the scene, but cannot approximate the spatially dense outputs required by segmentation or detection.

Delta-based sparse inference. A second line of work including RRM [20], CBinfer [2], Skip-Convolution [9], and DeltaCNN [22] maintains a pixel-level reference cache and propagates only the *difference* (delta) between the current frame and the reference through the network, computing only at affected output locations and reusing cached values elsewhere. Among these, DeltaCNN represents the most complete realization: it achieves end-to-end sparse propagation with truncation buffers that prevent error accumulation, enabling unbounded-length sequences without dense resets and pushing the efficiency of static-camera settings to its practical limit. MotionDeltaCNN [21] extends DeltaCNN to moving cameras by warping the cache with a single global homography before

computing the delta. However, a global transformation cannot capture locally heterogeneous motion—independently moving objects, depth-induced parallax, or mixed ego-motion and object motion—leaving large residual deltas that erode sparsity.

Shared limitation. Across both categories, the cache is treated as an *indivisible, scene-level entity*: it is either globally valid, globally stale, or aligned uniformly using a single transformation. Real-world mobile video, however, exhibits motion that is local and heterogeneous: a camera pan shifts the background uniformly while foreground objects move independently, and depth discontinuities create parallax that no single warp can reconcile. MotionDeltaCNN [21] reported that their homography alignment succeeds on only a small fraction of the DAVIS [23] sequences (14 out of 80) evaluated. The result is a *granularity mismatch*: the cache granularity is the whole scene, but motion granularity is per-region. This mismatch forces existing methods to either waste computation re-deriving content that has merely shifted, or sacrifice correctness by reusing misaligned features.

Note that several other techniques also reduce the cost of video analytics, including ROI filtering [3], [17], input resolution adaptation [14], [5], frame sampling that skips inference on selected frames [3], [17], and model compression via quantization or pruning [12], [11]. These strategies are orthogonal to feature cache reuse: they govern *which* frames, regions, or model to run, whereas caching determines *whether* to recompute or reuse at each spatial location within a given inference. FluxShard can be composed with any of them; this work focuses exclusively on the cache reuse axis.

REFERENCES

- [1] Hassan J. Al Dawasari, Muhammad Bilal, Muhammad Moinuddin, Kamran Arshad, and Khaled Assaleh. Deepvision: Enhanced drone detection and recognition in visible imagery through deep learning networks. *Sensors*, 23(21), 2023.
- [2] Lukas Cavigelli, Philippe Degen, and Luca Benini. CBinfer: Change-based inference for convolutional neural networks on video data.
- [3] Tiffany Yu-Han Chen, Lenin Ravindranath, Shuo Deng, Paramvir Bahl, and Hari Balakrishnan. Glimpse: Continuous, real-time object recognition on mobile devices. In *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*, SenSys ’15, pages 155–168. Association for Computing Machinery.
- [4] Hyomin Choi and Ivan V. Bajić. Scalable image coding for humans and machines. 31:2739–2754.
- [5] Kuntai Du, Qizheng Zhang, Anton Arapin, Haodong Wang, Zhengxu Xia, and Junchen Jiang. AccMPEG: Optimizing video encoding for video analytics.
- [6] Lingyu Duan, Jiaying Liu, Wenhan Yang, Tiejun Huang, and Wen Gao. Video coding for machines: A paradigm of collaborative compression and intelligent analytics. 29:8680–8695.
- [7] Leonardo Galteri, Lorenzo Seidenari, Marco Bertini, and Alberto Del Bimbo. Deep generative adversarial compression artifact removal. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 4836–4845. ISSN: 2380-7504.
- [8] Luyao Gao, Jianchun Liu, Hongli Xu, Sun Xu, Qianpiao Ma, and Liusheng Huang. Accelerating end-cloud collaborative inference via near bubble-free pipeline optimization. *arXiv preprint arXiv:2501.12388*, 2024.
- [9] Amirhossein Habibi, Davide Abati, Taco S. Cohen, and Babak Ehteshami Bejnordi. Skip-convolutions for efficient video processing. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2694–2703. ISSN: 2575-7075.
- [10] Beining Han, Meenal Parakh, Derek Geng, Jack A. Defay, Gan Luyang, and Jia Deng. FetchBench: A simulation benchmark for robot fetching.

- [11] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. In *Proceedings of the 29th International Conference on Neural Information Processing Systems - Volume 1*, volume 1 of *NIPS'15*, pages 1135–1143. MIT Press.
- [12] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2704–2713. IEEE.
- [13] Fatemeh Jalali, Morteza Khademi, Abbas Ebrahimi Moghadam, and Hadi Sadoghi Yazdi. Robust scene aware multi-object tracking for surveillance videos. 638:130114.
- [14] Junchen Jiang, Ganesh Ananthanarayanan, Peter Bodik, Siddhartha Sen, and Ion Stoica. Chameleon: scalable adaptation of video analytics. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication*, SIGCOMM '18, pages 253–266. Association for Computing Machinery.
- [15] Rahima Khanam and Muhammad Hussain. Yolov11: An overview of the key architectural enhancements. *arXiv preprint arXiv:2410.17725*, 2024.
- [16] Stefanos Laskaridis, Stylianos I Venieris, Mario Almeida, Ilias Leontiadis, and Nicholas D Lane. Spinn: synergistic progressive inference of neural networks over device and cloud. In *Proceedings of the 26th annual international conference on mobile computing and networking*, pages 1–15, 2020.
- [17] Yuanqi Li, Arthi Padmanabhan, Pengzhan Zhao, Yufei Wang, Guoqing Harry Xu, and Ravi Netravali. Reducto: On-camera filtering for resource-efficient real-time video analytics. In *Proceedings of the Annual conference of the ACM Special Interest Group on Data Communication on the applications, technologies, architectures, and protocols for computer communication*, SIGCOMM '20, pages 359–376. Association for Computing Machinery.
- [18] Mostafa Mahmoud, Kevin Siu, and Andreas Moshovos. Diffy: a déjà vu-free differential deep neural network accelerator. In *Proceedings of the 51st Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO-51, pages 134–147. IEEE Press.
- [19] Kien Nguyen, Feng Liu, Clinton Fookes, Sridha Sridharan, Xiaoming Liu, and Arun Ross. Person recognition in aerial surveillance: A decade survey. pages 1–1.
- [20] Bowen Pan, Wuwei Lin, Xiaolin Fang, Chaoqin Huang, Bolei Zhou, and Cewu Lu. Recurrent residual module for fast inference in videos. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1536–1545. ISSN: 2575-7075.
- [21] Mathias Parger, Chengcheng Tang, Thomas Neff, Christopher D. Twigg, Cem Keskin, Robert Wang, and Markus Steinberger. MotionDeltaCNN: Sparse CNN inference of frame differences in moving camera videos with spherical buffers and padded convolutions. pages 17292–17301.
- [22] Mathias Parger, Chengcheng Tang, Christopher D Twigg, Cem Keskin, Robert Wang, and Markus Steinberger. Deltacnn: End-to-end cnn inference of sparse frame differences in videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12497–12506, 2022.
- [23] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- [24] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *European Conference on Computer Vision (ECCV)*, sep 2018.
- [25] Qi Wu, Zipeng Fu, Xuxin Cheng, Xiaolong Wang, and Chelsea Finn. Helpful DoggyBot: Open-world object fetching using legged robots and vision-language models.
- [26] Jun Zhang, Mina Henein, Robert Mahony, and Viorela Ila. VDO-SLAM: A visual dynamic object-aware SLAM system.

A. Biographies and Author Photos

```
\begin{IEEEbiographynophoto}{Jane Doe}
Biography text here without a photo.
\end{IEEEbiographynophoto}
```

or a biography with a photo

```
\begin{IEEEbiography}[{\includegraphics
[width=1in,height=1.25in,clip,
```

```
keepaspectratio]{fig1.png}}]
{IEEE Publications Technology Team}
In this paragraph you can place
your educational, professional background
and research and other interests.
\end{IEEEbiography}
```

Please see the end of this document to see the output of these coding examples.

Jane Doe Biography text here without a photo.



IEEE Publications Technology Team In this paragraph you can place your educational, professional background and research and other interests.