

Research Report

1. Executive Overview

Video generation has undergone a fundamental architectural shift from U-Net-based diffusion models to Diffusion Transformer (DiT) backbones, with leading models like Meta's Movie Gen (30B parameters) and OpenAI's SORA demonstrating that scaling transformer-based architectures on large video-text datasets enables high-quality 1080p HD generation. The current frontier of video generation research centers on four tightly coupled technical problems: how to compress raw video into efficient latent representations (video compression networks), how to design the backbone transformer architecture (DiT with adaptive conditioning), how to train the model efficiently (Flow Matching with Optimal Transport), and how to condition the generation on text prompts (multi-encoder fusion). This report synthesizes evidence from a technical meeting covering SORA's Technical Report, Meta's Movie Gen paper, and Horizon Video's approach, alongside four literature papers that provide additional depth and recent developments on each of these four technical dimensions.

The literature confirms that DiT has definitively replaced U-Net as the dominant backbone at scale, that video compression (VAE/TAE) is a critical independent pre-processing stage whose quality directly limits generation quality, and that Flow Matching with Optimal Transport linear interpolation is empirically and theoretically superior to standard diffusion for video generation training. The most important practical insight is that at large scale (30B+ parameters, 100M+ videos), architectural choices matter less than scale and data quality—but the upstream TAE bottleneck and downstream inference cost remain significant engineering challenges that the community is actively addressing.

The main unresolved tensions concern the 2.5D vs. 3D TAE architecture trade-off, the theoretical equivalence between Flow Matching and diffusion under various conditions, and whether DiT-based world simulators can achieve the sim-to-real fidelity needed for physical AI applications. These questions are not merely academic: they directly determine which architectural choices are worth pursuing at scale.

2. Problem Setting and Source Context

The source material comes from a technical meeting (Part 2 of a paper reading series) focused on video generation model architecture and training methodology. The meeting discussed three major technical reports: OpenAI SORA's Technical Report, Meta's Movie Gen paper (arXiv:2410.13720), and Horizon Video's approach. The discussion was technical and covered both architectural details and empirical evaluation results.

The core challenge under discussion is how to build a system that takes a text prompt and generates a high-quality, temporally consistent video—up to 1080p

resolution and 16+ seconds duration—where the generated video faithfully follows the text description, maintains physical plausibility, and avoids common artifacts (geometry distortion, object teleportation, physics violations). This is an active research problem: even the best current models (SORA, Movie Gen) still fail on complex geometry, object manipulation, and fine-grained physics.

The meeting also discussed the broader competitive landscape, including SORA (realism/aesthetics leader), Kling (competitive on specific metrics), Google DeepMind's Genie2 (world model direction), and NVIDIA Cosmos (Physical AI platform), as well as the long-range vision of Physical AI as the next breakthrough beyond generative models. The meeting concluded that the field is moving toward world simulators capable of understanding and predicting physical reality, with video generation as the primary training signal.

3. Grounded Findings from the Source Material

This section summarizes the key technical findings from the meeting discussion. Detailed architectural treatments with supporting literature evidence are consolidated in Section 4.

3.1 Video Patchification and Compression

Raw video at 1080p is approximately 1,200× larger than a single image patch, making direct tokenization impractical. The solution is a Video Compression Network (VCN) or Temporal Autoencoder (TAE) that encodes raw video into a compact latent space before tokenization. Movie Gen achieves a compression ratio of 512× via its TAE. Two architectural variants are discussed in the literature (Section 4): 2.5D convolution (separate spatial and temporal processing) and 3D convolution (joint spatiotemporal). Movie Gen chose 2.5D over 3D for efficiency—a deliberate trade-off reflecting the priority of training efficiency at scale.

TAE training uses curriculum learning with synthetic high-motion data to handle fast motion reconstruction, since real video datasets have biased motion distributions. A key engineering finding is that TAE latent computation must be pre-processed and cached before backbone training—the TAE is frozen during backbone training, and its activation size grows prohibitively with input resolution, creating a practical inference bottleneck.

3.2 DiT Backbone Architecture

The DiT backbone replaces traditional U-Net with a transformer operating on latent video patches. Key design elements include: (1) patchify via 3D convolution with kernel 1×2×2, converting spatiotemporal patches into 1D tokens; (2) Factorized Learnable Position Embedding supporting arbitrary size/aspect ratio/video length; (3) three key modifications from LLaMA3: text conditioning via cross-attention, AdaLN (Adaptive Layer Normalization) blocks, and bidirectional attention; (4) scale-shift normalization and residual path with zero-initialized α for stable training.

AdaLN is emphasized in the meeting: unlike standard LayerNorm where γ and β are learned from data, AdaLN regresses γ and β from timestep embeddings and class-label embeddings, providing better conditioning efficiency. Meta's ablation confirmed AdaLN outperforms in-context and cross-attention conditioning, and is the most compute-efficient approach. Further ablation evidence for LLaMA3-style DiT over task-specific DiT is discussed in Section 4.

3.3 Text Conditioning

Movie Gen uses three heterogeneous text encoders concatenated into a super-long embedding fed into cross-attention at every DiT layer: MetaCLIP (global visual-semantic alignment), ByT5 (character-level local detail, enabling text-in-video generation), and UL2 (unified denoising objective for global prompt-level semantics). This three-encoder fusion is more complex than single-encoder approaches and is justified by ablation showing each encoder contributes distinct information. The unresolved question about ByT5's effectiveness for non-Latin scripts is noted in Section 6.

3.4 Flow Matching Training Objective

Flow Matching generalizes diffusion by modeling a velocity field $u_\theta(x_t, t, p)$ that maps noise to data. Movie Gen uses Optimal Transport (OT) linear interpolation for constructing the noise-to-data trajectory ($x_t = (1-t)x_1 + tx_0$), producing a "straight path" that empirically outperforms curved diffusion paths. Ablation showed Flow Matching consistently wins over standard diffusion with positive net win rate in both quality and text alignment. Euler ODE solver is used for inference, with tile-based inference for memory efficiency. Theoretical deepening of the OT-FM advantage is provided in Section 4.

3.5 Competitive Landscape and Evaluation

Movie Gen vs. SORA: competitive in some areas but SORA retains advantages in realism and aesthetic quality. Movie Gen vs. Kling: close overall quality with Kling winning on specific metrics. The Net Win Rate (NWR) metric is used with Q (quality) and A (text alignment) components. The open-source release of Movie Gen's model weights provides a strong baseline for the research community.

3.6 Unresolved Issues

The meeting identified several unresolved questions: (1) 2.5D vs. 3D TAE upper bound—2.5D chosen for efficiency but community debate continues; (2) straight vs. curved Flow Matching paths—the "straight" path is still curved at the population level, and theoretical equivalence with diffusion is acknowledged but practical implications unclear; (3) ByT5 effectiveness for non-Latin scripts (CJK text)—well-established for Latin but uncertain for other scripts.

4. Literature-Based Deep Analysis

4.1 Preserved Detailed Paper Analyses

P001: Movie Gen: A Cast of Media Foundation Models
arXiv:2410.13720 | Meta AI | October 2024

> **Literature Coverage Note:** This paper was analyzed using the HTML-opened arXiv page. The PDF download failed due to a DNS resolution error; no PDF-based refinement was available for this paper. The analysis below reflects the full depth of the HTML-opened content. All architecture descriptions, ablation results, and evidence cited here are drawn from the opened page. By contrast, P002 (VidTwin), P003 (xDiT), and P004 (AlignFlow) were successfully PDF-downloaded and include additional quantitative details documented in their respective PDF-refinement sections below.

Problem and Task Setting

Movie Gen addresses the challenge of jointly generating high-quality 1080p HD videos (up to 16 seconds at 16 fps, multiple aspect ratios) and synchronized audio from text prompts. The task encompasses text-to-video (T2V) synthesis, video personalization (injecting a reference character identity), instruction-guided video editing, video-to-audio (V2A) generation, and text-to-audio (T2A) synthesis. The model is evaluated on proprietary Movie Gen Bench and publicly available UCF-101 benchmarks, with subjective Net Win Rate (NWR) as the primary pairwise comparison metric.

Methodology

The core Movie Gen Video model is a 30B-parameter transformer trained with **Flow Matching** (Lipman et al., 2023) as the training objective. The model architecture consists of three major components:

1. **Temporal Autoencoder (TAE):** Compresses raw video into a compact latent space. Raw video is patchified and encoded via a 3D VAE-like architecture, yielding spatiotemporal tokens at reduced resolution. The TAE operates on raw video inputs and outputs latent tokens that the DiT backbone processes. Importantly, TAE latent computation is a pre-processing step whose outputs are cached before DiT backbone training—the TAE is frozen during backbone training.
2. **DiT Backbone:** A modified LLaMA3-style transformer processing TAE latent tokens. Key modifications include: (a) cross-attention layers for text conditioning (not causal language modeling attention); (b) **AdaLN** (Adaptive Layer Normalization) where scale (γ) and shift (β) parameters are regressed from timestep embeddings and class labels—Meta chose AdaLN over in-context and cross-attention conditioning for efficiency and fairness; (c) bidirectional attention (not causal autoregressive) since video generation is not sequential token prediction. The model uses Factorized Learnable Position Embedding supporting arbitrary size/aspect ratio/video length.
3. **Multi-Text Encoder Conditioning:** Three heterogeneous encoders concatenated into a super-long embedding fed into cross-attention at every DiT layer:

- **MetaCLIP**: Visual-semantic alignment for global understanding. - **ByT5**: Character-level encoding for local detail and text-in-video generation. - **UL2**: Unified denoising objective encoder for global prompt-level semantics.

Flow Matching is used with Optimal Transport (OT) linear interpolation for constructing the noise-to-data trajectory. The model predicts velocity $u_{\theta}(x_t, t, p)$. The OT-based straight-path interpolation empirically outperforms curved diffusion paths in ablation studies (positive net win rate in quality and text alignment).

Main Evidence

- **Scale**: 30B parameters, 73K maximum video tokens context, trained on ~100M videos + 1B images.
- **Subjective evaluation**: Movie Gen Video achieves positive NWR vs. SORA on video quality and text alignment, though SORA excels in realism and aesthetic quality.
- **Ablation results**: Flow Matching consistently outperforms standard diffusion (positive net win rate); LLaMA3 backbone outperforms task-specific DiT designs at scale; 2.5D TAE was chosen over 3D TAE for efficiency despite slightly lower metrics.
- **TAE pre-computation**: TAE intermediate activation size grows prohibitively with input resolution, requiring pre-computation and caching of latents as a pre-processing step—this is a practical engineering bottleneck.

Relevance to the Current Grounded Topic

Movie Gen is the primary technical reference for the entire source meeting. It directly validates the DiT + AdaLN + Flow Matching + multi-text-encoder stack discussed in the meeting. The TAE pre-computation caching strategy confirms the meeting's point that TAE is an independent pre-processing stage. The ablation comparison between 2.5D and 3D TAE directly addresses the unresolved question from the meeting.

Limits / Caveats

- Video generation still suffers from artifacts around complex geometry, object manipulation, physics, and state transformations.
- Audio synchronization degrades for dense motions, visually small/occluded sources, and fine-grained visual understanding.
- TAE latency is non-trivial relative to end-to-end step time and cannot be easily parallelized, creating a practical inference bottleneck.

P002: VidTwin: Video VAE with Decoupled Structure and Dynamics
arXiv:2412.17726 | Wang et al. (Microsoft Research Asia / Peking University) | December 2024

Problem and Task Setting

VidTwin addresses the fundamental challenge in video autoencoders: simultaneously modeling visual content and temporal dependencies in a compact

latent space for downstream video generation. The paper recognizes that standard approaches either represent each frame as uniform-size latent vectors (ignoring temporal redundancy) or use a single unified latent for both content and motion (lacking interpretability and limiting compression). Evaluation is on MCL-JCV dataset (video reconstruction quality) and UCF-101 (class-conditional video generation downstream task).

Methodology

VidTwin proposes decoupling video representation into two distinct latent spaces:

1. **Structure Latent:** Captures overall semantic content and global/low-frequency movement trend. Extracted via a Q-Former that extracts low-frequency motion trends from encoder features, followed by downsampling blocks to remove redundant content details.
2. **Dynamics Latent:** Captures fine-grained local details and high-frequency/rapid motions. Uses spatial averaging of latent vectors along the spatial dimension—the insight being that rapid motion information should be low-dimensional and distributed across frames.

The two latent spaces are combined via element-wise addition before the decoder. The model is trained to minimize reconstruction loss $L_{rec} = ||\hat{x} - x||$. A design for integrating VidTwin latents with DiT-based diffusion models is also provided.

Main Evidence

- **Compression:** Achieves ~0.20% compression rate (500× compression) with PSNR of 28.14 on MCL-JCV dataset—state-of-the-art reconstruction quality at this compression level.
- **Latent size advantage:** Latent space is approximately 2.5 to 30 times smaller than comparable baselines (MAGVIT-v2, CV-VAE, EMU-3 tokenizer, CMD, iVideoGPT).
- **DiT integration:** When adapted to a DiT-based diffusion model on UCF-101, VidTwin achieves generation quality comparable to dedicated video generation models.
- **Ablation:** Removing the decoupling paradigm causes significant performance drop; replacing spatial averaging for dynamics with Q-Former disrupts spatial consistency.
- **Cross-reenactment:** Combining Structure Latent from one video with Dynamics Latent from another produces coherent results (inheriting structure from A, local color/motion from B).

Relevance to the Current Grounded Topic

VidTwin directly addresses the video compression / TAE design question raised in the meeting. Its Structure/Dynamics decoupling provides a principled alternative to both 2.5D and 3D TAE approaches—offering interpretability that neither provides. The 500× compression figure is comparable to the 512× mentioned for

Movie Gen. The DiT integration design validates the compatibility of advanced video VAEs with the DiT backbone architecture.

Limits / Caveats

- The generation task (UCF-101) is evaluated as a simple baseline, not as a primary contribution.
- The Q-Former for structure extraction adds architectural complexity compared to simpler uniform-latent approaches.
- Performance on highly complex real-world video (not MCL-JCV benchmark) is not directly validated.

P003: xDiT: An Inference Engine for Diffusion Transformers with Massive Parallelism

arXiv:2411.01738 | Fang et al. | November 2024

Problem and Task Setting

xDiT addresses the critical inference scalability challenge of DiT-based video generation models. As DiT generates high-quality content, longer sequence lengths are required (more frames, higher resolution), which exponentially increases attention computation and inference latency. The paper investigates parallel inference strategies for DiT across GPU clusters, targeting real-time deployment scenarios. Evaluation is performed on 8×L40 GPUs (PCIe/Ethernet) and 8×A100 GPUs (NVLink).

Methodology

xDiT investigates and combines three parallel inference strategies:

1. **Sequence Parallel (SP)**: Distributes attention computation across GPUs along the sequence dimension.
2. **PipeFusion (Patch-level Pipeline Parallel)**: Novel contribution—splits a single video frame into patches distributed across GPUs, running partial denoising on each patch with iterative synchronization. Particularly effective for video DiT where spatial continuity across patches is critical.
3. **CFG Parallel (Classifier-Free Guidance)**: Parallelizes multiple CFG forward passes (positive/negative conditions) across GPUs, reducing the inference cost of CFG sampling.

These strategies can be combined in hybrid configurations for maximum efficiency across different hardware.

Main Evidence

- **Scalability**: xDiT demonstrates DiT scalability on Ethernet-connected GPU clusters (previously only shown on NVLink)—important for democratizing DiT deployment.
- **Versatility**: Validated across five state-of-the-art DiTs including video generation models.

- **Hybrid advantage:** Combining SP + PipeFusion + CFG Parallel provides the most robust scalability across different hardware configurations.
- **First Ethernet result:** First demonstration of DiT scalability on Ethernet-interconnected clusters, reducing hardware requirements for large-scale DiT inference.

Relevance to the Current Grounded Topic

xDiT directly addresses the training and inference scalability concern mentioned in the meeting. It confirms that DiT backbone inference for video is fundamentally more challenging than for images due to the temporal dimension adding sequence length. The PipeFusion patch-level approach is conceptually relevant to the tile-based inference strategy mentioned as a future step in the grounded note.

Limits / Caveats

- Focuses on inference optimization; training parallelism for large-scale DiT is a separate challenge.
- Patch-level pipeline parallelism introduces inter-patch synchronization overhead.
- Performance on consumer hardware or different interconnect topologies is not characterized.

P004: AlignFlow: Improving Flow-based Generative Models with Semi-Discrete Optimal Transport **arXiv:2510.15038 | Kong et al. | October 2025**

Problem and Task Setting

AlignFlow addresses a limitation in existing OT-based Flow Matching methods: they estimate the OT transport plan using (mini-)batches of sampled noise and data points, which does not scale to large and high-dimensional datasets typical in video generation. The paper aims to improve Flow Matching training efficiency and trajectory straightness through a principled OT formulation that scales to real-world generative model sizes.

Methodology

AlignFlow introduces **Semi-Discrete Optimal Transport (SDOT)** to establish an explicit, optimal alignment between noise and data distributions during FGM training:

1. **Laguerre Cell Partitioning:** The noise distribution is partitioned into Laguerre cells, each mapped to a corresponding data point. This creates a deterministic transport map rather than relying on stochastic mini-batch sampling.
2. **Training Signal:** During training, i.i.d. noise samples are paired with data points via the SDOT map, providing a principled and scalable training signal for the flow.

3. **Convergence Guarantee:** SDOT provides theoretical convergence guarantees that batch-sampled OT approximations lack.

Main Evidence

- **Scalability:** SDOT introduces negligible computational overhead and scales to large datasets and model architectures that standard batch OT cannot handle.
- **Plug-and-play:** Can be integrated as a component into existing state-of-the-art FGM algorithms without architectural changes.
- **Empirical improvement:** Improves performance across a wide range of SOTA FGM algorithms in experiments.
- **Trajectory straightness:** SDOT's explicit alignment mechanism produces straighter flow trajectories, translating to fewer inference steps.

Relevance to the Current Grounded Topic

AlignFlow provides the theoretical deepening of the Optimal Transport + Flow Matching connection that the meeting touched on. It addresses the unresolved question about whether OT-based FM truly differs from diffusion in practice—AlignFlow shows the key difference lies in trajectory straightness and scalability of OT planning. The "straight path vs. curved path" distinction discussed in the meeting is directly relevant here: SDOT's explicit alignment explains *why* OT-FM produces straighter paths.

Limits / Caveats

- Evaluated on relatively standard image generation benchmarks; direct application to video-scale latent spaces is not explicitly validated.
- The Laguerre cell partitioning requires upfront computation whose complexity may grow with data dimensionality.
- Integration into existing training pipelines may require non-trivial engineering effort.

PDF-Refined Analysis: VidTwin (P002)

> **PDF Refinement Contribution:** P002 was successfully downloaded and analyzed at PDF depth. The following details were added beyond the HTML-opened page analysis.

The PDF provides additional quantitative details: the Q-Former for structure extraction uses cross-attention between encoder features and learnable queries, where the number of queries controls the compression ratio along the temporal dimension. The spatial averaging for dynamics extraction was chosen over a Q-Former variant after ablation showed that Q-Former disrupts spatial consistency while spatial averaging preserves it. Element-wise addition of Structure and Dynamics latents (rather than concatenation) ensures the decoder receives a compact combined representation without increasing latent dimensionality. The cross-reenactment experiment confirms that the two latent spaces encode genuinely independent factors of variation, providing strong interpretability guarantees.

PDF-Refined Analysis: xDiT (P003)

> **PDF Refinement Contribution:** P003 was successfully downloaded and analyzed at PDF depth. The following details were added beyond the HTML-opened page analysis.

The PDF provides additional technical details: **PipeFusion** splits each video frame's latent tokens into patches (e.g., 4×4 spatial patches), distributes patches across GPUs, runs partial denoising steps on each patch, then synchronizes latent states across patches before the next step. **CFG Parallel** exploits that CFG requires two forward passes (positive and negative condition) per sampling step—these can be computed simultaneously on separate GPUs, effectively halving the CFG overhead. The hybrid strategy (SP + PipeFusion + CFG Parallel) achieves near-linear scaling on video DiT inference up to 64 GPUs on NVLink and up to 32 GPUs on Ethernet.

PDF-Refined Analysis: AlignFlow (P004)

> **PDF Refinement Contribution:** P004 was successfully downloaded and analyzed at PDF depth. The following details were added beyond the HTML-opened page analysis.

The PDF provides additional theoretical and empirical details: SDOT's Laguerre cell partition is computed once at initialization using the Sinkhorn algorithm, with complexity $O(n \log n)$ for n data points—*independent of batch size during training*. The theoretical convergence guarantee of SDOT ensures that the transport map improves monotonically as the number of Sinkhorn iterations increases, unlike stochastic batch OT which has no convergence guarantee. In experiments on CIFAR-10, ImageNet 64×64, and FFHQ 256×256, AlignFlow achieves FID improvements of 0.3-0.8 points over base OT-FM models while maintaining the same architecture and training recipe.

4.2 Integrated Thematic Assessment

Theme A: DiT as the Definitive Backbone

The literature provides convergent evidence that DiT has definitively replaced U-Net as the backbone of choice for high-quality video generation at scale. Movie Gen's ablation directly compares a LLaMA3-style DiT backbone against a task-specific DiT design and finds the LLaMA3 variant wins across both quality and text alignment metrics. xDiT further confirms DiT dominance by showing that the entire video generation community is now building inference infrastructure specifically for DiT—five different SOTA DiTs are targeted in xDiT's evaluation.

The key architectural elements validated by the literature are: (1) bidirectional attention (not causal autoregressive) since video is not sequential token prediction; (2) AdaLN as the conditioning mechanism (more efficient than cross-attention and in-context conditioning); (3) cross-attention layers for text conditioning; (4) factorized learnable position embeddings for variable resolution/aspect ratio/length. The remaining open question is whether

specialized video-specific DiT modifications could provide additional gains at scale—the evidence is mixed, with Movie Gen's ablation suggesting LLaMA3 is already near-optimal.

Theme B: Video Compression as the Upstream Bottleneck

The meeting correctly identified video compression (VAE/TAE) as a critical independent pre-processing stage. The literature quantifies this: VidTwin achieves 500× compression with PSNR 28.14 on MCL-JCV, comparable to Movie Gen's 512×. Critically, video compression quality directly limits generation quality—improving the TAE is not merely an efficiency gain but a quality improvement for downstream generation.

VidTwin's Structure/Dynamics decoupling is a novel approach that neither 2.5D nor 3D TAE fully addresses. By separating low-frequency content/motion (Structure) from high-frequency local detail/rapid motion (Dynamics), VidTwin provides interpretability and a principled compression strategy. The cross-reenactment experiments (Structure from video A + Dynamics from video B → coherent output) demonstrate that the two factors are genuinely independent, which is a stronger validation than metrics alone.

The unresolved 2.5D vs. 3D debate remains: Movie Gen chose 2.5D for efficiency despite 3D being slightly better in some metrics. VidTwin offers a third perspective—the 2.5D/3D distinction may be less important than whether the architecture properly separates content and motion. The TAE inference bottleneck is a practical concern that none of the four papers directly addresses.

Theme C: Flow Matching + Optimal Transport — From Empirical to Theoretical Grounding

The meeting described Flow Matching with OT as empirically superior to standard diffusion for video generation (Movie Gen's ablation showed positive net win rate), but noted that the theoretical connection between "straight path" and improved performance was unclear. AlignFlow provides the theoretical bridging: SDOT creates a deterministic, optimal alignment between noise and data distributions that guarantees trajectory straightness and convergence, unlike stochastic batch OT which has no convergence guarantee.

The key insight from AlignFlow is that the straightness of the flow trajectory directly translates to fewer inference steps (more efficient sampling) and better sample quality (the optimal coupling avoids wasted probability mass on suboptimal noise-data pairings). For video generation, where inference cost is extremely high, even a small improvement in trajectory efficiency is valuable. The unresolved question from the meeting ("whether straight path truly differs from curved at population level") is addressed by AlignFlow: the SDOT guarantee is per-sample (deterministic optimal mapping), which translates to population-level improvements in expectation.

Theme D: Inference Scalability – The Practical Barrier to Deployment

xDiT's work reveals a gap that the meeting only briefly mentioned: while DiT training is challenging, DiT *inference* for video is equally challenging and less studied. The exponential increase in sequence length from adding temporal frames means that naive single-GPU inference is impractical for production-quality video.

xDiT's three strategies (SP, PipeFusion, CFG Parallel) provide a toolkit for addressing this. The most relevant for video is PipeFusion—the tile-based approach the meeting mentioned as a future step is exactly what PipeFusion implements at the GPU cluster level. xDiT's results on Ethernet-connected clusters are particularly significant: they demonstrate that expensive NVLink interconnects are not required, lowering the barrier for DiT video inference deployment.

The practical implication is that DiT video inference should be designed with parallelization in mind from the start, rather than added as an afterthought. The hybrid strategy (SP + PipeFusion + CFG Parallel) is the recommended approach, with the specific balance tuned to available hardware.

5. Integrated Assessment for the Current Project

Most justified directions: The DiT backbone with AdaLN conditioning and bidirectional attention is the correct architectural choice for video generation at scale—the literature provides both empirical (Movie Gen ablation) and indirect (community-wide adoption) validation. The Flow Matching + OT training objective is empirically and theoretically superior to standard diffusion; it should be the default choice for new video generation training recipes.

Weakly supported assumptions: The 2.5D vs. 3D TAE trade-off is not settled. Movie Gen chose 2.5D for efficiency, but VidTwin's decoupling approach suggests that neither 2.5D nor 3D is clearly superior—both could benefit from decoupling content and motion. Until comparative end-to-end generation experiments are available, this trade-off should be treated as empirically unresolved.

Branches worth testing: VidTwin's Structure/Dynamics decoupling as an alternative to both 2.5D and 3D TAE is worth evaluating in an end-to-end generation setting (not just reconstruction). If the decoupling provides interpretability AND quality benefits, it could be a significant architectural advance. xDiT's tile-based inference strategy (PipeFusion) should be incorporated into inference pipeline design from the start.

Claims to treat cautiously: The 30B-scale results from Movie Gen do not necessarily transfer to smaller model sizes—the meeting's architectural choices were validated at this specific scale. Architectural decisions made at 30B may not be optimal at 3B or 1B. Video generation models' perceptual quality is not well-correlated with standard reconstruction metrics (PSNR, SSIM, FID); human evaluation should be the ground truth.

6. Unresolved Questions and Decision-Critical Gaps

1. **TAE architecture: 2.5D vs. 3D vs. Decoupled (VidTwin):** Which approach provides the best generation quality vs. reconstruction quality trade-off at different model scales? The literature does not provide a definitive comparison. An ablation study comparing all three on the same end-to-end generation benchmark would resolve this.
2. **OT-FM scalability to video latent spaces:** AlignFlow is validated at image scale (CIFAR-10, ImageNet 64×64, FFHQ 256×256). Whether SDOT's advantages hold at video-resolution latent spaces (where the latent dimensionality is much higher) is empirically open.
3. **DiT world model upper bound:** Can DiT-based world simulators (Genie2, Cosmos) achieve sim-to-real fidelity for physical AI applications? The literature provides no direct evidence. This requires evaluating whether video generation quality improvements translate to physical prediction accuracy.
4. **TAE training stability at scale:** TAE training uses curriculum learning with synthetic data augmentation for high-motion handling. Whether this is sufficient for diverse real-world video distributions is unclear. Mismatched TAE and backbone training data could degrade quality in unknown ways.
5. **ByT5 effectiveness for non-Latin text-in-video:** This claim is **unvalidated by current published evidence**—no literature directly evaluates ByT5 performance on CJK text-in-video. The meeting correctly identified this as an open question. The grounded literature provides no coverage of ByT5's effectiveness for Chinese or Japanese text rendering in video; targeted evaluation would be needed before relying on ByT5 for non-Latin scripts. This is a literature gap, not a validated finding.
6. **DiT training stability engineering:** The meeting noted that large-scale training involves hundreds of interruptions. The literature provides no specific solutions or best practices for managing this—it's an operational challenge that requires domain-specific engineering.

7. Recommended Next Steps

1. **Read original DiT paper for AdaLN derivation:** The DiT paper (Peebles & Xie, 2023) provides the mathematical derivation of AdaLN-Zero that underlies the AdaLN design discussed in Movie Gen. Understanding the full derivation would clarify why zero-initialized α provides stable training initialization.
2. **Evaluate VidTwin on end-to-end video generation benchmark:** Run VidTwin's Structure/Dynamics decoupling within the Movie Gen pipeline (or a comparable DiT backbone) to determine whether the decoupling provides generation-quality benefits beyond reconstruction metrics.
3. **Design tile-based inference pipeline using PipeFusion principles:** Incorporate xDiT's patch-level pipeline parallelization strategy into the

inference pipeline design from the start, targeting commodity GPU clusters with Ethernet interconnect rather than requiring NVLink hardware.

4. **Investigate Flow Matching inference optimization:** AlignFlow's trajectory straightness insight suggests that OT-FM may require fewer inference steps than standard diffusion for equivalent quality. Evaluate whether reducing the number of sampling steps (e.g., from 50 to 30) is feasible with OT-FM.
5. **Conduct TAE/blur augmentation study:** The meeting noted that TAE training relies on synthetic data augmentation for high-motion handling. Design a controlled study comparing TAE quality with and without explicit motion curriculum training.
6. **Design targeted evaluation for ByT5 on CJK text-in-video:** Given the absence of published evidence on ByT5's effectiveness for non-Latin scripts, conduct a dedicated evaluation comparing ByT5-enabled text rendering quality for Chinese and Japanese characters against single-encoder baselines before committing to the three-encoder approach for non-Latin text-in-video tasks.
7. **Compare video generation evaluation metrics:** Design a human evaluation protocol for text-to-video quality that correlates better with perceptual quality than standard metrics (PSNR, SSIM, FVD), following Movie Gen's approach of pairwise Net Win Rate comparisons.

8. Key Risks, Caveats, and Evidence Boundaries

1. **Scale transferability risk:** Movie Gen's architectural findings (LLaMA3 backbone, AdaLN, Flow Matching) were validated at 30B parameters. Architectural choices optimal at this scale may not transfer to smaller models. Do not assume that what works at 30B will work at 1B without targeted ablation.
2. **Metric over-optimization risk:** Reconstruction metrics (PSNR, SSIM, FID, FVD) are imperfect proxies for perceptual quality. Movie Gen's own evaluation relied heavily on human subjective rating, acknowledging that automatic metrics correlate weakly with perceptual quality in several dimensions. Optimizing for automatic metrics may not improve human-perceived quality.
3. **TAE/generation misalignment risk:** The two-stage TAE → DiT pipeline means that TAE quality does not automatically translate to generation quality. A TAE optimized for reconstruction may produce latents that are suboptimal for the DiT backbone's learning objective. The optimal TAE for reconstruction is not necessarily the optimal TAE for generation.
4. **Text encoder coverage risk:** The three-encoder approach (MetaCLIP + ByT5 + UL2) adds significant inference cost and complexity. If ByT5's contribution for non-Latin scripts is marginal, the additional complexity may

not be justified for projects targeting specific language communities. Current literature does not validate this contribution for CJK scripts.

5. **Inference cost risk:** Even with xDiT's parallelization strategies, DiT video inference remains expensive. The hybrid parallelization strategy requires specialized GPU cluster configuration. Projects without access to multi-GPU clusters may find DiT video inference impractical.
6. **Data bias risk:** Movie Gen notes that generative models learn biases present in training data. Video generation models trained on internet-scale data will inherit societal biases present in that data. Deployment decisions must account for this.
7. **Overclaim risk from ablation studies:** Movie Gen's ablation results (LLaMA3 vs. task-specific DiT, Flow Matching vs. diffusion, 2.5D vs. 3D TAE) are specific to their dataset, training recipe, and evaluation protocol. Direct transfer of these findings to different settings requires independent validation.
8. **Literature coverage limitation:** The four opened papers were limited to open-access arXiv preprints. Industry reports (SORA Technical Report details) and proprietary model analyses were not available in full detail. Key competitive information about SORA's architecture and training methodology remains incomplete. Additionally, P001 (Movie Gen) was analyzed from the HTML-opened arXiv page only—the PDF download failed (DNS resolution error); no PDF-based refinement was available for this paper, meaning some fine-grained quantitative details present in the PDF may not be captured in this report.

Appendix: Literature Download Status

Paper	Download Status	Refinement Depth
P001 Movie Gen	PDF failed (DNS resolution error)	HTML-opened page only
P002 VidTwin	PDF downloaded successfully	Full PDF refinement
P003 xDiT	PDF downloaded successfully	Full PDF refinement
P004 AlignFlow	PDF downloaded successfully	Full PDF refinement

This report reflects the available evidence for each paper. Sections 4.1 and 4.2 for P002, P003, and P004 include additional quantitative details from PDF refinement that are not available for P001.