

To evaluate whether Fourier Transform (FT) could be applied to feature reuse in the context of the diffusion transformer (DiT) acceleration methods described in the provided documents (TeaCache, FORA, ToCa, and TaylorSeer), let's analyze the potential applicability of FT based on the characteristics of these methods and the properties of FT. The discussion will focus on how FT might enhance or complement feature reuse strategies in DiT models for image and video generation, considering the temporal and spatial nature of the diffusion process.

---

## Understanding Feature Reuse in DiT Models

The four methods (TeaCache, FORA, ToCa, and TaylorSeer) accelerate DiT models by leveraging the similarity or predictability of features across timesteps or tokens during the iterative denoising process. Here's a brief recap of their approaches:

- **TeaCache:** Estimates output differences using timestep-embedding modulated noisy inputs and selectively caches outputs to skip redundant computations.
- **FORA:** Caches and reuses intermediate outputs from attention and MLP layers, exploiting feature similarity across timesteps.
- **ToCa:** Selects specific tokens for caching based on temporal redundancy and error propagation, optimizing at a fine-grained token level.
- **TaylorSeer:** Predicts future features using Taylor series expansion, modeling feature trajectories across timesteps to enable high-ratio acceleration.

These methods rely on the temporal redundancy or continuity of features across denoising steps, where features (or their differences) evolve gradually. The challenge lies in efficiently identifying, storing, or predicting these features to reduce computational overhead while maintaining generation quality.

Fourier Transform, which decomposes a signal into its frequency components, is widely used in signal processing, image analysis, and other domains to analyze patterns in the frequency domain. It could potentially be applied to feature reuse if the features or their differences exhibit periodic or frequency-based patterns across timesteps or spatial dimensions.

---

## Potential Applications of Fourier Transform in Feature Reuse

To determine whether FT can be used, let's consider how it might fit into the feature reuse paradigm, focusing on the temporal and spatial aspects of DiT features.

### 1. Temporal Analysis of Features Across Timesteps

The diffusion process in DiT models involves iterative denoising over multiple timesteps, where features (e.g., intermediate outputs or token representations) evolve smoothly or exhibit redundancy. This is evident in:

- **TaylorSeer**: Features form stable trajectories across timesteps, with derivatives showing consistent patterns (as seen in PCA projections).
- **FORA and ToCa**: High similarity between features in consecutive timesteps, particularly in later denoising stages.
- **TeaCache**: Differences between outputs vary across timesteps, with input differences correlating to output differences.

How FT Could Apply:

- **Frequency Decomposition of Feature Trajectories:**

- Features across timesteps can be treated as a time-series signal for each token or layer. Applying FT could transform these features into the frequency domain, revealing dominant frequency components that characterize their evolution.
- Low-frequency components might capture the smooth, gradual changes in features (as observed in TaylorSeer's stable trajectories), while high-frequency components could represent noise or rapid changes that are less suitable for caching.
- By filtering out high-frequency components (e.g., using a low-pass filter), one could isolate the stable, redundant parts of the feature signal for caching or prediction, potentially improving the efficiency of methods like FORA or ToCa.

- **Compression for Caching:**

- FT could compress feature representations by retaining only significant frequency components, reducing memory overhead for caching. This is particularly relevant for FORA, which caches attention and MLP layer outputs, and ToCa, which caches token-specific features.
- For example, instead of storing full feature tensors, one could store their Fourier coefficients (or a subset thereof), reconstructing them when needed. This could reduce storage costs and accelerate retrieval.

- **Prediction Enhancement for TaylorSeer:**

- TaylorSeer predicts future features using Taylor series based on derivatives. FT could complement this by analyzing the frequency content of feature derivatives (as shown in TaylorSeer's Figure 1(b)). If derivatives exhibit periodic patterns, FT could help model these patterns, potentially simplifying the prediction process or improving accuracy for distant timesteps.
- For instance, FT could identify periodic trends in feature changes, allowing a hybrid approach where low-frequency trends are predicted using FT-based methods and higher-order derivatives are handled by Taylor series.

## Challenges:

- The denoising process is not strictly periodic, so FT might not capture all relevant patterns. The smooth trajectories observed in TaylorSeer suggest that low-frequency components dominate, but non-periodic changes could complicate FT-based analysis.
- Computing FT for high-dimensional feature tensors (e.g., token embeddings or layer outputs) could introduce computational overhead, potentially offsetting the acceleration gains unless optimized (e.g., using Fast Fourier Transform, FFT).

## 2. Spatial Analysis of Features Across Tokens

ToCa highlights that different tokens exhibit varying temporal redundancy and error propagation, suggesting spatial heterogeneity in feature behavior. Features in DiT models are often spatially structured, especially for image and video generation, where tokens correspond to patches or regions in the input.

**How FT Could Apply:**

- **Spatial Redundancy in Tokens:**

- FT is commonly used in image processing to analyze spatial frequencies (e.g., low frequencies for smooth regions, high frequencies for edges). Applying 2D FT to feature maps (e.g., token representations in a spatial grid) could identify spatially redundant regions suitable for caching.
- For ToCa, FT could enhance token selection by identifying tokens with low-frequency spatial content, which are likely to be more redundant and less sensitive to caching errors. This could complement ToCa's temporal redundancy and error propagation scores.

- **Feature Compression:**

- Similar to temporal compression, 2D FT could compress spatial feature maps by retaining only low-frequency components, reducing the memory footprint for caching in methods like FORA or ToCa.
- For video generation, 3D FT (spatial-temporal) could capture both spatial and temporal redundancies, potentially identifying stable regions across frames for caching.

- **Error Propagation Analysis:**

- ToCa considers error propagation through attention layers. FT could analyze the spatial frequency content of attention maps to identify tokens with minimal influence on others, refining ToCa's selection criteria.

### **Challenges:**

- The spatial structure of DiT features (e.g., patch-based tokens) may not align perfectly with traditional image-based FT, requiring adaptation to handle token grids.
- High-dimensional feature spaces in DiT models could make FT computationally expensive, especially for large models like OpenSora or FLUX.

## **3. Enhancing Difference Estimation in TeaCache**

TeaCache estimates output differences using timestep-embedding modulated noisy inputs and refines them with polynomial fitting. FT could potentially enhance this process:

- **Frequency-Based Difference Analysis:**

- By applying FT to the sequence of noisy inputs or timestep embeddings, TeaCache could analyze the frequency components of input differences. Low-frequency differences might indicate stable timesteps suitable for caching, while high-frequency differences could signal timesteps requiring full computation.
- This could provide an alternative or complementary method to polynomial fitting, potentially simplifying the scaling factor estimation.

- **Noise Characterization:**

- The noisy inputs in diffusion models contain structured noise that evolves over timesteps. FT could characterize this noise in the frequency domain, helping TeaCache identify timesteps where input differences are dominated by low-frequency components (indicating high correlation with output differences).

#### **Challenges:**

- The correlation between input and output differences (as exploited by TeaCache) may not translate directly to the frequency domain, requiring careful validation.
  - FT might introduce complexity that undermines TeaCache's goal of low computational overhead.
- 

## **Practical Considerations and Feasibility**

While FT shows promise for enhancing feature reuse, several practical aspects need consideration:

- **Computational Cost:** FFT algorithms are efficient ( $O(n \log n)$ ), but applying FT to high-dimensional feature tensors across multiple timesteps or tokens could add overhead. This is a concern for training-free methods like TeaCache, FORA, ToCa, and TaylorSeer, which prioritize minimal computation.
  - **Feature Dimensionality:** DiT features are high-dimensional (e.g., token embeddings or layer outputs), requiring multidimensional FT (e.g., 2D for spatial, 3D for spatial-temporal). This increases complexity compared to 1D FT used in traditional signal processing.
  - **Model Integration:** Integrating FT into existing caching pipelines would require modifying the feature processing pipeline, potentially affecting the plug-and-play nature of these methods.
  - **Quality Trade-offs:** FT-based compression or prediction might introduce artifacts (e.g., due to lossy frequency truncation), which could degrade generation quality. Careful tuning of frequency thresholds would be needed to match the near-lossless performance of methods like ToCa or TaylorSeer.
- 

## Specific Scenarios Where FT Could Be Beneficial

Based on the methods' characteristics, here are specific scenarios where FT could enhance feature reuse:

1. **TaylorSeer**: FT could model periodic or semi-periodic patterns in feature derivatives, complementing Taylor series predictions. For example, FT could identify dominant frequencies in feature trajectories, allowing more accurate long-range predictions for high-acceleration ratios (e.g., 5×).
  2. **ToCa**: FT could refine token selection by analyzing spatial frequency content, prioritizing tokens with low-frequency features for caching. This could reduce error propagation in attention layers by focusing on stable, low-frequency regions.
  3. **FORA**: FT-based compression could reduce the memory footprint of cached attention and MLP layer outputs, enabling faster retrieval and storage, especially for large-scale models like PIXART- $\alpha$ .
  4. **TeaCache**: FT could enhance input difference estimation by analyzing the frequency content of noisy inputs or timestep embeddings, potentially simplifying the polynomial fitting step.
- 

## Limitations and Challenges

- **Non-Periodic Nature**: The diffusion process is iterative and not inherently periodic, which may limit FT's effectiveness compared to methods like Taylor series (used in TaylorSeer) that directly model continuous trajectories.
  - **High-Dimensional Data**: Applying FT to high-dimensional feature tensors requires significant memory and computation, potentially negating acceleration benefits unless optimized (e.g., using sparse FT or dimensionality reduction).
  - **Validation Needs**: The effectiveness of FT in capturing feature redundancy or predicting future features needs empirical validation. The documents (e.g., TaylorSeer's PCA results) suggest smooth trajectories, but it's unclear if frequency-based patterns are prominent enough to justify FT.
  - **Integration Complexity**: Modifying existing methods to incorporate FT could disrupt their training-free nature, requiring additional preprocessing or post-processing steps.
- 

## Conclusion

Fourier Transform could theoretically be applied to feature reuse in DiT models, particularly for:

- **Compressing cached features** to reduce memory overhead (benefiting FORA and ToCa).
- **Analyzing temporal or spatial redundancy** to enhance token or timestep selection (benefiting ToCa and TeaCache).
- **Modeling feature trajectories** to improve long-range predictions (benefiting TaylorSeer).

However, its practical utility is limited by the non-periodic nature of the diffusion process, the high dimensionality of features, and the computational cost of FT in this context. FT is most likely to be effective as a complementary tool for:

- **Low-frequency feature extraction** to identify stable, cacheable components.
- **Compression of feature representations** to optimize storage in caching mechanisms.

To integrate FT effectively, one would need to:

1. Use FFT for computational efficiency.
2. Focus on low-frequency components to capture smooth feature changes.
3. Validate empirically using datasets like ImageNet or video generation benchmarks (as used in the documents) to ensure quality preservation.

Compared to the existing methods, FT might not replace core strategies like TaylorSeer's predictive approach or ToCa's token-wise selection but could enhance them in specific scenarios, particularly for compression or redundancy analysis. Further experimentation would be needed to confirm its benefits, especially given the success of simpler, non-frequency-based methods like polynomial fitting (TeaCache) or derivative-based prediction (TaylorSeer).

If you'd like, I can explore a specific implementation (e.g., applying FT to ToCa's token selection) or simulate a small-scale FT-based analysis using a toy dataset to illustrate its potential. Let me know!