

# Understanding the CacheView Class: Key-Value Cache Management for Transformer Models

Code Analysis with Numerical Examples

June 11, 2025

## 1 Introduction

The CacheView class is a sophisticated cache management system designed for efficient key-value storage and retrieval in transformer-based models. This system is particularly crucial for optimizing inference performance by avoiding redundant computations of attention keys and values that have already been processed.

## 2 Class Overview

The CacheView class manages four main components:

- `cache_k`: Cached key tensors
- `cache_v`: Cached value tensors
- `kv_seqlens`: Sequence lengths for each batch element
- `metadata`: Additional metadata for cache operations

## 3 Detailed Method Analysis

### 3.1 Constructor (`__init__`)

```
1 def __init__(  
2     self,  
3     cache_k: torch.Tensor,  
4     cache_v: torch.Tensor,  
5     metadata: CacheInputMetadata,  
6     kv_seqlens: torch.Tensor,  
7 ):  
8     self.cache_k = cache_k  
9     self.cache_v = cache_v  
10    self.kv_seqlens = kv_seqlens  
11    self.metadata = metadata
```

The constructor initializes the cache with pre-allocated tensors and metadata.

**Numerical Example:** Let's assume we have a batch of 2 sequences with the following dimensions:

$$\text{cache\_k} \in \mathbb{R}^{2 \times 512 \times 8 \times 64} \quad (1)$$

$$\text{cache\_v} \in \mathbb{R}^{2 \times 512 \times 8 \times 64} \quad (2)$$

$$\text{kv\_seqlens} = [128, 256] \quad (3)$$

Where:

- Batch size = 2
- Maximum sequence length = 512
- Number of KV heads = 8
- Head dimension = 64
- Current sequence lengths are 128 and 256 tokens respectively

## 3.2 Update Method

```

1 def update(self, xk: torch.Tensor, xv: torch.Tensor) -> None:
2     """
3         to_cache_mask masks the last [max_seq_len] tokens in each sequence
4     """
5     n_kv_heads, head_dim = self.cache_k.shape[-2:]
6     flat_cache_k = self.cache_k.view(-1, n_kv_heads, head_dim)
7     flat_cache_v = self.cache_v.view(-1, n_kv_heads, head_dim)
8
9     flat_cache_k.index_copy_(0, self.metadata.cache_positions,
10                             xk[self.metadata.to_cache_mask])
11    flat_cache_v.index_copy_(0, self.metadata.cache_positions,
12                             xv[self.metadata.to_cache_mask])

```

This method updates the cache with new key-value pairs using efficient indexing operations.

### Step-by-Step Analysis:

#### 1. Extract dimensions:

$$n\_kv\_heads = 8, \quad head\_dim = 64$$

#### 2. Flatten cache tensors:

$$\text{flat\_cache\_k} : \mathbb{R}^{2 \times 512 \times 8 \times 64} \rightarrow \mathbb{R}^{1024 \times 8 \times 64} \quad (4)$$

$$\text{flat\_cache\_v} : \mathbb{R}^{2 \times 512 \times 8 \times 64} \rightarrow \mathbb{R}^{1024 \times 8 \times 64} \quad (5)$$

**3. Index-based copying:** The `index_copy_` operation efficiently places new key-value pairs at specific positions in the flattened cache.

**Numerical Example:** Suppose we have new tokens to cache:

$$xk \in \mathbb{R}^{10 \times 8 \times 64} \text{ (10 new tokens)} \quad (6)$$

$$\text{to\_cache\_mask} = [\text{True}, \text{True}, \text{False}, \text{True}, \dots] \text{ (mask for 5 tokens)} \quad (7)$$

$$\text{cache\_positions} = [128, 129, 131, 132, 133] \text{ (positions in flat cache)} \quad (8)$$

The operation `flat_cache_k.index_copy_(0, cache_positions, xk[to_cache_mask])` will:

- Take the masked tokens from `xk` (5 tokens)
- Place them at positions [128, 129, 131, 132, 133] in the flattened cache

### 3.3 Interleave KV Method

```

1 def interleave_kv(self, xk: torch.Tensor, xv: torch.Tensor) -> Tuple[torch
2 .Tensor, torch.Tensor]:
3     """
4     This is a naive implementation and not optimized for speed.
5     """
6     assert xk.ndim == xv.ndim == 3 # (B * T, H, D)
7     assert xk.shape == xv.shape
8
9     if all([s == 0 for s in self.metadata.seqlens]):
10         # No cache to interleave
11         return xk, xv
12
13     # Make it a list of [(T, H, D)]
14     xk: Tuple[torch.Tensor] = torch.split(xk, self.metadata.seqlens)
15     xv: Tuple[torch.Tensor] = torch.split(xv, self.metadata.seqlens)
16     assert len(xk) == len(self.kv_seqlens), f"Batch size is {len(self.
17     kv_seqlens)}, got {len(xk)}"
18
19     # Order elements in cache by position by unrotating
20     cache_k = [unrotate(t, s) for t, s in zip(self.cache_k, self.
21     kv_seqlens)]
22     cache_v = [unrotate(t, s) for t, s in zip(self.cache_v, self.
23     kv_seqlens)]
24
25     interleaved_k = interleave_list(cache_k, list(xk))
26     interleaved_v = interleave_list(cache_v, list(xv))
27
28     return torch.cat(interleaved_k, dim=0), torch.cat(interleaved_v, dim
29 =0)

```

This method combines cached key-value pairs with new ones, handling sequence-level operations.

#### Detailed Process:

1. **Input Validation:** Ensures input tensors are 3-dimensional with shape  $(B \times T, H, D)$

## 2. Sequence Splitting:

$$xk \rightarrow [(T_1, H, D), (T_2, H, D), \dots] \quad (9)$$

$$xv \rightarrow [(T_1, H, D), (T_2, H, D), \dots] \quad (10)$$

**3. Cache Unrotation:** The `unrotate` function reorders cached elements by their original positions.

**4. Interleaving:** Combines cached and new key-value pairs in the correct sequential order.

**Numerical Example:** Consider a batch with 2 sequences:

$$\text{Input: } xk \in \mathbb{R}^{6 \times 8 \times 64} \text{ (total 6 new tokens)} \quad (11)$$

$$\text{seqlens} = [3, 3] \text{ (3 tokens per sequence)} \quad (12)$$

$$\text{kv\_seqlens} = [128, 256] \text{ (cached lengths)} \quad (13)$$

After splitting:

$$xk[0] \in \mathbb{R}^{3 \times 8 \times 64} \text{ (new tokens for seq 1)} \quad (14)$$

$$xk[1] \in \mathbb{R}^{3 \times 8 \times 64} \text{ (new tokens for seq 2)} \quad (15)$$

After interleaving:

$$\text{Final sequence 1: } \mathbb{R}^{(128+3) \times 8 \times 64} = \mathbb{R}^{131 \times 8 \times 64} \quad (16)$$

$$\text{Final sequence 2: } \mathbb{R}^{(256+3) \times 8 \times 64} = \mathbb{R}^{259 \times 8 \times 64} \quad (17)$$

## 3.4 Properties

The class provides several useful properties:

```

1 @property
2 def max_seq_len(self) -> int:
3     return self.cache_k.shape[1]
4
5 @property
6 def key(self) -> torch.Tensor:
7     return self.cache_k[: len(self.kv_seqlens)]
8
9 @property
10 def value(self) -> torch.Tensor:
11     return self.cache_v[: len(self.kv_seqlens)]
12
13 @property
14 def prefill(self) -> bool:
15     return self.metadata.prefill
16
17 @property
18 def mask(self) -> AttentionBias:
19     return self.metadata.mask

```

**Property Explanations:**

- `max_seq_len`: Returns the maximum sequence length the cache can handle
- `key/value`: Returns only the active portion of the cache (up to current batch size)
- `prefill`: Indicates whether this is a prefill operation
- `mask`: Returns the attention bias/mask for the current operation

## 4 Use Case and Benefits

### 4.1 Efficiency Gains

The CacheView class provides several efficiency benefits:

1. **Memory Reuse**: Avoids recomputing attention keys and values for previously processed tokens
2. **Batch Processing**: Handles multiple sequences simultaneously
3. **Flexible Updates**: Supports incremental updates without full recomputation

### 4.2 Typical Usage Pattern

1. **Initialization**: Create cache with pre-allocated tensors
2. **Update**: Add new key-value pairs as tokens are processed
3. **Interleave**: Combine cached and new data for attention computation
4. **Access**: Use properties to get current cache state

## 5 Mathematical Representation

Let's represent the cache operations mathematically:

### 5.1 Cache Update Operation

Given new keys  $K_{new} \in \mathbb{R}^{T_{new} \times H \times D}$  and cache positions  $P \in \mathbb{N}^{T_{cache}}$ :

$$C_k[P[i]] \leftarrow K_{new}[M[i]] \quad \forall i \text{ where } M[i] = \text{True}$$

Where  $M$  is the `to_cache_mask` and  $C_k$  is the flattened cache.

### 5.2 Interleaving Operation

For sequence  $s$  with cached length  $L_s$  and new length  $T_s$ :

$$K_{final}^{(s)} = \text{concat}(\text{unrotate}(C_k^{(s)}, L_s), K_{new}^{(s)})$$

The final concatenated result across all sequences:

$$K_{batch} = \text{concat}(K_{final}^{(1)}, K_{final}^{(2)}, \dots, K_{final}^{(B)})$$

## 6 Conclusion

The `CacheView` class represents a sophisticated approach to managing key-value caches in transformer models. By providing efficient update mechanisms, flexible interleaving operations, and clean property access, it enables significant performance improvements in inference scenarios where maintaining state across multiple forward passes is crucial.

The mathematical operations, while complex, are designed to minimize computational overhead while maintaining the semantic correctness of attention computations. This makes it particularly valuable for applications requiring real-time or near-real-time transformer inference.

## 7 Comprehensive Real-World Example: Chatbot Inference

Let's walk through a complete real-world scenario where a chatbot is processing multiple conversations simultaneously. This example will demonstrate every aspect of the `CacheView` class with concrete numbers.

### 7.1 Scenario Setup

**Context:** A customer service chatbot is handling 3 simultaneous conversations:

- **Conversation 1:** Technical support (currently 15 tokens processed)
- **Conversation 2:** Product inquiry (currently 8 tokens processed)
- **Conversation 3:** Billing question (currently 22 tokens processed)

**Model Configuration:**

$$\text{Batch size (B)} = 3 \quad (18)$$

$$\text{Max sequence length} = 128 \quad (19)$$

$$\text{Number of KV heads (H)} = 4 \quad (20)$$

$$\text{Head dimension (D)} = 32 \quad (21)$$

$$\text{New tokens per step} = 1 \text{ (autoregressive generation)} \quad (22)$$

### 7.2 Initial Cache State

The cache tensors are initialized as:

$$\text{cache\_k} \in \mathbb{R}^{3 \times 128 \times 4 \times 32} \quad (23)$$

$$\text{cache\_v} \in \mathbb{R}^{3 \times 128 \times 4 \times 32} \quad (24)$$

$$\text{kv\_seqlens} = [15, 8, 22] \quad (25)$$

The cache structure looks like this:

Sequence	Used Positions	Free Positions	Total
Conv 1	0-14 (15 tokens)	15-127	128
Conv 2	0-7 (8 tokens)	8-127	128
Conv 3	0-21 (22 tokens)	22-127	128

### 7.3 Step 1: Processing New User Messages

New user inputs arrive:

- **Conv 1:** "What" (token ID: 2841)
- **Conv 2:** "How" (token ID: 1374)
- **Conv 3:** "When" (token ID: 4829)

**Input tensors:**

$$xk \in \mathbb{R}^{3 \times 4 \times 32} \text{ (3 new key vectors)} \quad (26)$$

$$xv \in \mathbb{R}^{3 \times 4 \times 32} \text{ (3 new value vectors)} \quad (27)$$

**Metadata configuration:**

$$\text{to\_cache\_mask} = [True, True, True] \quad (28)$$

$$\text{cache\_positions} = [15, 8, 22] \text{ (next available positions)} \quad (29)$$

$$\text{seqlens} = [1, 1, 1] \text{ (one new token each)} \quad (30)$$

### 7.4 Step 2: Cache Update Operation

The update method processes the new tokens:

**Dimension extraction:**

$$n\_kv\_heads = 4, \quad head\_dim = 32$$

**Cache flattening:**

$$\text{flat\_cache\_k} : \mathbb{R}^{3 \times 128 \times 4 \times 32} \rightarrow \mathbb{R}^{384 \times 4 \times 32} \quad (31)$$

$$\text{flat\_cache\_v} : \mathbb{R}^{3 \times 128 \times 4 \times 32} \rightarrow \mathbb{R}^{384 \times 4 \times 32} \quad (32)$$

**Position mapping in flattened cache:**

$$\text{Conv 1, pos 15} \rightarrow \text{flat index } 0 \times 128 + 15 = 15 \quad (33)$$

$$\text{Conv 2, pos 8} \rightarrow \text{flat index } 1 \times 128 + 8 = 136 \quad (34)$$

$$\text{Conv 3, pos 22} \rightarrow \text{flat index } 2 \times 128 + 22 = 278 \quad (35)$$

**Index copy operation:**

```

1 # Pseudocode for the actual operation
2 flat_cache_k[15] = xk[0]      # Conv 1's new key
3 flat_cache_k[136] = xk[1]     # Conv 2's new key
4 flat_cache_k[278] = xk[2]     # Conv 3's new key
5
6 flat_cache_v[15] = xv[0]      # Conv 1's new value
7 flat_cache_v[136] = xv[1]     # Conv 2's new value
8 flat_cache_v[278] = xv[2]     # Conv 3's new value

```

Updated sequence lengths:

$$\text{kv\_seqlens} = [16, 9, 23]$$

## 7.5 Step 3: Generating Next Tokens

The model now generates responses. For the attention computation, we need to combine cached keys/values with any new computation keys/values.

New computation inputs (for attention):

$$xk_{attn} \in \mathbb{R}^{3 \times 4 \times 32} \text{ (current step's keys)} \quad (36)$$

$$xv_{attn} \in \mathbb{R}^{3 \times 4 \times 32} \text{ (current step's values)} \quad (37)$$

## 7.6 Step 4: Interleaving Operation

The `interleave_kv` method combines cached and new data:

**Input validation:**

$$xk_{attn}.shape = (3, 4, 32) \text{ 3D tensor} \quad (38)$$

$$xv_{attn}.shape = (3, 4, 32) \text{ matching shapes} \quad (39)$$

**Sequence splitting:** Since each sequence contributes 1 token:

$$xk_{split} = [(1, 4, 32), (1, 4, 32), (1, 4, 32)] \quad (40)$$

$$xv_{split} = [(1, 4, 32), (1, 4, 32), (1, 4, 32)] \quad (41)$$

**Cache extraction and unrotation:**

$$\text{cache\_k}[0] \in \mathbb{R}^{128 \times 4 \times 32} \text{ (Conv 1 cache)} \quad (42)$$

$$\text{cache\_k}[1] \in \mathbb{R}^{128 \times 4 \times 32} \text{ (Conv 2 cache)} \quad (43)$$

$$\text{cache\_k}[2] \in \mathbb{R}^{128 \times 4 \times 32} \text{ (Conv 3 cache)} \quad (44)$$

After `unrotate` with sequence lengths [16, 9, 23]:

$$\text{unrotated\_cache\_k}[0] \in \mathbb{R}^{16 \times 4 \times 32} \text{ (active part)} \quad (45)$$

$$\text{unrotated\_cache\_k}[1] \in \mathbb{R}^{9 \times 4 \times 32} \text{ (active part)} \quad (46)$$

$$\text{unrotated\_cache\_k}[2] \in \mathbb{R}^{23 \times 4 \times 32} \text{ (active part)} \quad (47)$$

**Interleaving process:** For each conversation, combine cached tokens with new token:

$$\text{final\_k}[0] = \text{concat}(\text{unrotated\_cache\_k}[0], xk_{split}[0]) \quad (48)$$

$$\in \mathbb{R}^{(16+1) \times 4 \times 32} = \mathbb{R}^{17 \times 4 \times 32} \quad (49)$$

$$\text{final\_k}[1] = \text{concat}(\text{unrotated\_cache\_k}[1], xk_{split}[1]) \quad (50)$$

$$\in \mathbb{R}^{(9+1) \times 4 \times 32} = \mathbb{R}^{10 \times 4 \times 32} \quad (51)$$

$$\text{final\_k}[2] = \text{concat}(\text{unrotated\_cache\_k}[2], xk_{split}[2]) \quad (52)$$

$$\in \mathbb{R}^{(23+1) \times 4 \times 32} = \mathbb{R}^{24 \times 4 \times 32} \quad (53)$$

**Final concatenation:**

$$K_{final} = \text{concat}(\text{final\_k}[0], \text{final\_k}[1], \text{final\_k}[2]) \quad (54)$$

$$\in \mathbb{R}^{(17+10+24) \times 4 \times 32} = \mathbb{R}^{51 \times 4 \times 32} \quad (55)$$

## 7.7 Step 5: Attention Computation

The model now has complete key-value pairs for attention:

- **Conv 1:** 17 tokens (16 cached + 1 new)
- **Conv 2:** 10 tokens (9 cached + 1 new)
- **Conv 3:** 24 tokens (23 cached + 1 new)

The attention computation uses these complete sequences while maintaining conversation boundaries through appropriate masking.

## 7.8 Step 6: Next Iteration

The model generates responses:

- **Conv 1:** "about" (token ID: 1234)
- **Conv 2:** "much" (token ID: 5678)
- **Conv 3:** "will" (token ID: 9012)

These new tokens will be cached at positions [16, 9, 23] respectively, and the process repeats.

**Updated state after this iteration:**

$$\text{kv_seqlens} = [17, 10, 24] \quad (56)$$

$$\text{cache\_positions} = [17, 10, 24] \text{ (for next update)} \quad (57)$$

## 7.9 Performance Analysis

Memory efficiency:

- Without caching: Would need to recompute  $17 + 10 + 24 = 51$  token embeddings
- With caching: Only compute 3 new token embeddings
- Speedup:**  $\frac{51}{3} = 17 \times$  fewer computations

Memory usage:

$$\text{Cache memory} = 3 \times 128 \times 4 \times 32 \times 2 (\text{k}+\text{v}) \times 4 \text{ bytes (float32)} \quad (58)$$

$$= 196,608 \text{ bytes} = 192 \text{ KB per conversation batch} \quad (59)$$

This demonstrates how the `CacheView` class enables efficient real-time conversation handling by avoiding redundant computations while maintaining perfect attention semantics.

## 8 Ultra-Detailed Example: Complete Tensor Operations Walkthrough

To fully understand the inner workings of the `CacheView` class, let's examine a comprehensive example that shows actual tensor values, complete mathematical operations, and multiple cache update cycles.

### 8.1 Initial Setup with Concrete Tensor Values

**Scenario:** A small-scale transformer processing 2 sequences with the following specifications:

$$\text{Batch size} = 2 \quad (60)$$

$$\text{Max sequence length} = 8 \quad (61)$$

$$\text{Number of KV heads} = 2 \quad (62)$$

$$\text{Head dimension} = 4 \quad (63)$$

$$\text{Initial cached lengths} = [3, 2] \quad (64)$$

Initial cache state:

```
1 # cache_k shape: (2, 8, 2, 4)
2 # cache_v shape: (2, 8, 2, 4)
3 # kv_seqlens: [3, 2]
4
5 cache_k = [
6     [ # Sequence 0
7         [[1.0, 2.0, 3.0, 4.0], [5.0, 6.0, 7.0, 8.0]], # Position 0
8         [[1.1, 2.1, 3.1, 4.1], [5.1, 6.1, 7.1, 8.1]], # Position 1
9         [[1.2, 2.2, 3.2, 4.2], [5.2, 6.2, 7.2, 8.2]], # Position 2
```

```

10     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 3 (
11     unused)
12     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 4 (
13     unused)
14     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 5 (
15     unused)
16     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 6 (
17     unused)
18     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]] # Position 7 (
19     unused)
20 ],
21 [ # Sequence 1
22     [[2.0, 3.0, 4.0, 5.0], [6.0, 7.0, 8.0, 9.0]], # Position 0
23     [[2.1, 3.1, 4.1, 5.1], [6.1, 7.1, 8.1, 9.1]], # Position 1
24     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 2 (
25     unused)
26     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 3 (
27     unused)
28     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 4 (
29     unused)
30     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 5 (
31     unused)
32     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]], # Position 6 (
33     unused)
34     [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]] # Position 7 (
35     unused)
36 ]
37 ]

```

## 8.2 Iteration 1: Adding New Tokens

New input tokens:

- Sequence 0: 2 new tokens
- Sequence 1: 1 new token

Input tensors:

```

1 # xk shape: (3, 2, 4) - total 3 new tokens across sequences
2 xk = [
3     [[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]], # Seq 0, Token 1
4     [[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]], # Seq 0, Token 2
5     [[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]] # Seq 1, Token 1
6 ]
7
8 # xv shape: (3, 2, 4) - corresponding value vectors
9 xv = [
10    [[0.9, 0.8, 0.7, 0.6], [0.5, 0.4, 0.3, 0.2]], # Seq 0, Token 1
11    [[0.91, 0.81, 0.71, 0.61], [0.51, 0.41, 0.31, 0.21]], # Seq 0, Token 2
12    [[0.8, 0.7, 0.6, 0.5], [0.4, 0.3, 0.2, 0.1]] # Seq 1, Token 1
13 ]

```

Metadata configuration:

```

1 metadata = CacheInputMetadata(
2     to_cache_mask = [True, True, True],      # Cache all 3 tokens
3     cache_positions = [3, 4, 2],           # Positions: [seq0_pos3,
4     seq0_pos4, seq1_pos2]
5     seqlens = [2, 1],                   # New tokens per sequence
6     # ... other metadata fields
)

```

## 8.3 Update Operation Step-by-Step

### Step 1: Extract dimensions

```

1 n_kv_heads, head_dim = cache_k.shape[-2:] # n_kv_heads=2, head_dim=4

```

### Step 2: Flatten cache tensors

```

1 # Original: (2, 8, 2, 4) -> Flattened: (16, 2, 4)
2 flat_cache_k = cache_k.view(-1, 2, 4)
3
4 # Mapping:
5 # flat_cache_k[0] = cache_k[0][0] (seq 0, pos 0)
6 # flat_cache_k[1] = cache_k[0][1] (seq 0, pos 1)
7 # ...
8 # flat_cache_k[8] = cache_k[1][0] (seq 1, pos 0)
9 # flat_cache_k[9] = cache_k[1][1] (seq 1, pos 1)
10 #

```

### Step 3: Index copy operations

```

1 # cache_positions = [3, 4, 2] maps to flat indices:
2 # Position 3 in seq 0 -> flat_index = 0*8 + 3 = 3
3 # Position 4 in seq 0 -> flat_index = 0*8 + 4 = 4
4 # Position 2 in seq 1 -> flat_index = 1*8 + 2 = 10
5
6 # Copy operations:
7 flat_cache_k[3] = xk[0] # [[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]]
8 flat_cache_k[4] = xk[1] # [[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]]
9 flat_cache_k[10] = xk[2] # [[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]]
10
11 flat_cache_v[3] = xv[0] # [[0.9, 0.8, 0.7, 0.6], [0.5, 0.4, 0.3, 0.2]]
12 flat_cache_v[4] = xv[1] # [[0.91, 0.81, 0.71, 0.61], [0.51, 0.41, 0.31,
13 0.21]]
13 flat_cache_v[10] = xv[2] # [[0.8, 0.7, 0.6, 0.5], [0.4, 0.3, 0.2, 0.1]]

```

### Updated cache state:

```

1 # After update, kv_seqlens = [5, 3]
2 # Sequence 0 now has 5 tokens (positions 0,1,2,3,4)
3 # Sequence 1 now has 3 tokens (positions 0,1,2)
4
5 updated_cache_k = [
6     [ # Sequence 0 - now 5 tokens
7         [[1.0, 2.0, 3.0, 4.0], [5.0, 6.0, 7.0, 8.0]],          # Position 0
8         (old)
9         [[1.1, 2.1, 3.1, 4.1], [5.1, 6.1, 7.1, 8.1]],          # Position 1
10        (old)
11    ]
12 ]

```

```

9      [[1.2, 2.2, 3.2, 4.2], [5.2, 6.2, 7.2, 8.2]],           # Position 2
10     ([[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]],          # Position 3
11     ([[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]],          # Position 4
12     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 5
13     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 6
14     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]])         # Position 7
15     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]]         # Position 8
16   ],
17   [ # Sequence 1 - now 3 tokens
18     [[2.0, 3.0, 4.0, 5.0], [6.0, 7.0, 8.0, 9.0]],          # Position 0
19     ([[2.1, 3.1, 4.1, 5.1], [6.1, 7.1, 8.1, 9.1]],          # Position 1
20     ([[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]],          # Position 2
21     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 3
22     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 4
23     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 5
24     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 6
25     ([[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]]         # Position 7
26   ]

```

## 8.4 Interleaving Operation Detailed Walkthrough

Now let's trace through the interleaving operation for attention computation:

**Input for interleaving:**

```

1 # New tokens for current attention computation
2 xk_attention = [
3   [[7.0, 6.0, 5.0, 4.0], [3.0, 2.0, 1.0, 0.0]],    # Seq 0 current token
4   [[6.0, 5.0, 4.0, 3.0], [2.0, 1.0, 0.0, -1.0]]    # Seq 1 current token
5 ]
6
7 # Metadata for current step
8 current_seqlens = [1, 1] # One new token per sequence

```

**Step 1: Input validation and splitting**

```

1 # xk_attention.shape = (2, 2, 4)      3D tensor
2 # Split by sequence lengths [1, 1]:
3 xk_split = [
4   [[[7.0, 6.0, 5.0, 4.0], [3.0, 2.0, 1.0, 0.0]]],    # Seq 0: (1, 2, 4)
5   [[[6.0, 5.0, 4.0, 3.0], [2.0, 1.0, 0.0, -1.0]]]    # Seq 1: (1, 2, 4)

```

6 ]

### Step 2: Extract and unrotate cache

```

1 # Extract active cache portions based on kv_seqlens = [5, 3]
2 cache_k_active = [
3     [ # Sequence 0: 5 tokens
4         [[1.0, 2.0, 3.0, 4.0], [5.0, 6.0, 7.0, 8.0]],      # Position 0
5         [[1.1, 2.1, 3.1, 4.1], [5.1, 6.1, 7.1, 8.1]],      # Position 1
6         [[1.2, 2.2, 3.2, 4.2], [5.2, 6.2, 7.2, 8.2]],      # Position 2
7         [[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]],      # Position 3
8         [[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]]       # Position 4
9     ],
10    [ # Sequence 1: 3 tokens
11        [[2.0, 3.0, 4.0, 5.0], [6.0, 7.0, 8.0, 9.0]],      # Position 0
12        [[2.1, 3.1, 4.1, 5.1], [6.1, 7.1, 8.1, 9.1]],      # Position 1
13        [[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]]        # Position 2
14    ]
15 ]
16
17 # unrotate() ensures proper positional ordering (assuming no rotation
18     needed here)
unrotated_cache_k = cache_k_active # Same in this example

```

### Step 3: Interleave cached and new tokens

```

1 # Combine cached tokens with new tokens for each sequence
2 interleaved_k = [
3     [ # Sequence 0: 5 cached + 1 new = 6 total tokens
4         [[1.0, 2.0, 3.0, 4.0], [5.0, 6.0, 7.0, 8.0]],      # Cached pos
5         0
6         [[1.1, 2.1, 3.1, 4.1], [5.1, 6.1, 7.1, 8.1]],      # Cached pos
7         1
8         [[1.2, 2.2, 3.2, 4.2], [5.2, 6.2, 7.2, 8.2]],      # Cached pos
9         2
10        [[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]],      # Cached pos
11        3
12        [[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]],      # Cached pos
13        4
14        [[7.0, 6.0, 5.0, 4.0], [3.0, 2.0, 1.0, 0.0]]        # NEW token
15    ],
16    [ # Sequence 1: 3 cached + 1 new = 4 total tokens
17        [[2.0, 3.0, 4.0, 5.0], [6.0, 7.0, 8.0, 9.0]],      # Cached pos
18        0
19        [[2.1, 3.1, 4.1, 5.1], [6.1, 7.1, 8.1, 9.1]],      # Cached pos
20        1
21        [[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]],      # Cached pos
22        2
23        [[6.0, 5.0, 4.0, 3.0], [2.0, 1.0, 0.0, -1.0]]        # NEW token
24    ]
25 ]

```

### Step 4: Final concatenation

```

1 # Concatenate all sequences for batch attention computation
2 # Final shape: (10, 2, 4) - 10 total tokens across both sequences
3 final_k = [

```

```

4 # Sequence 0 tokens (6 tokens):
5 [[1.0, 2.0, 3.0, 4.0], [5.0, 6.0, 7.0, 8.0]],
6 [[1.1, 2.1, 3.1, 4.1], [5.1, 6.1, 7.1, 8.1]],
7 [[1.2, 2.2, 3.2, 4.2], [5.2, 6.2, 7.2, 8.2]],
8 [[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]],
9 [[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]],
10 [[7.0, 6.0, 5.0, 4.0], [3.0, 2.0, 1.0, 0.0]],
11 # Sequence 1 tokens (4 tokens):
12 [[2.0, 3.0, 4.0, 5.0], [6.0, 7.0, 8.0, 9.0]],
13 [[2.1, 3.1, 4.1, 5.1], [6.1, 7.1, 8.1, 9.1]],
14 [[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]],
15 [[6.0, 5.0, 4.0, 3.0], [2.0, 1.0, 0.0, -1.0]]
16 ]

```

## 8.5 Iteration 2: Additional Cache Growth

Let's continue with another iteration to show cache evolution:

New input (iteration 2):

```

1 # Adding 1 token to sequence 0, 2 tokens to sequence 1
2 xk_iter2 = [
3     [[10.0, 11.0, 12.0, 13.0], [14.0, 15.0, 16.0, 17.0]], # Seq 0, Token
4         1
5     [[20.0, 21.0, 22.0, 23.0], [24.0, 25.0, 26.0, 27.0]], # Seq 1, Token
6         1
7     [[20.1, 21.1, 22.1, 23.1], [24.1, 25.1, 26.1, 27.1]]    # Seq 1, Token
8         2
9 ]
10
11 # Metadata for iteration 2
12 metadata_iter2 = CacheInputMetadata(
13     to_cache_mask = [True, True],           # Cache both new tokens
14     cache_positions = [5, 3],               # Next positions in cache
15     seqlens = [1, 2],                      # One token for seq0, 2 tokens for
16     seq1
17     prefill = False                       # Generation phase
18 )

```

Cache state after iteration 2:

```

1 # Updated kv_seqlens = [6, 5]
2 final_cache_k = [
3     [ # Sequence 0 - now 6 tokens (positions 0-5)
4         [[1.0, 2.0, 3.0, 4.0], [5.0, 6.0, 7.0, 8.0]],          # Position 0
5         [[1.1, 2.1, 3.1, 4.1], [5.1, 6.1, 7.1, 8.1]],          # Position 1
6         [[1.2, 2.2, 3.2, 4.2], [5.2, 6.2, 7.2, 8.2]],          # Position 2
7         [[9.0, 8.0, 7.0, 6.0], [5.0, 4.0, 3.0, 2.0]],          # Position 3
8         [[9.1, 8.1, 7.1, 6.1], [5.1, 4.1, 3.1, 2.1]],          # Position 4
9         [[10.0, 11.0, 12.0, 13.0], [14.0, 15.0, 16.0, 17.0]], # Position 5
10        (NEW)
11         [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],          # Position 6
12        (unused)
13         [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]]           # Position 7
14        (unused)

```

```

12  ],
13  [ # Sequence 1 - now 5 tokens (positions 0-4)
14  [[2.0, 3.0, 4.0, 5.0], [6.0, 7.0, 8.0, 9.0]],      # Position 0
15  [[2.1, 3.1, 4.1, 5.1], [6.1, 7.1, 8.1, 9.1]],      # Position 1
16  [[8.0, 7.0, 6.0, 5.0], [4.0, 3.0, 2.0, 1.0]],      # Position 2
17  [[20.0, 21.0, 22.0, 23.0], [24.0, 25.0, 26.0, 27.0]], # Position 3
18  (NEW)
19  [[20.1, 21.1, 22.1, 23.1], [24.1, 25.1, 26.1, 27.1]], # Position 4
20  (NEW)
21  [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],       # Position 5
22  (unused)
23  [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]],       # Position 6
24  (unused)
25  [[0.0, 0.0, 0.0, 0.0], [0.0, 0.0, 0.0, 0.0]]        # Position 7
26  (unused)
27  ]
28 ]

```

## 8.6 Performance Metrics for This Example

Cache efficiency analysis:

$$\text{Initial state: } 3 + 2 = 5 \text{ tokens cached} \quad (65)$$

$$\text{After iteration 1: } 5 + 3 = 8 \text{ tokens cached} \quad (66)$$

$$\text{After iteration 2: } 6 + 5 = 11 \text{ tokens cached} \quad (67)$$

$$\text{New computations: } 3 + 3 = 6 \text{ tokens total} \quad (68)$$

$$\text{Computation savings: } \frac{11}{6} = 83\% \text{ fewer calculations} \quad (69)$$

**Memory utilization:**

$$\text{Allocated memory: } 2 \times 8 \times 2 \times 4 \times 4 = 512 \text{ bytes} \quad (70)$$

$$\text{Used memory: } (6 + 5) \times 2 \times 4 \times 4 = 352 \text{ bytes} \quad (71)$$

$$\text{Utilization rate: } \frac{352}{512} = 68.75\% \quad (72)$$

This ultra-detailed walkthrough demonstrates the precise tensor operations, memory management, and computational efficiency achieved by the `CacheView` class. Every step shows concrete values, making the abstract cache operations tangible and verifiable.

## 9 Complete Transformer Forward Passes: 3-Pass Deep Dive

To understand how `CacheView` integrates into the complete transformer architecture, let's trace through 3 full forward passes from raw tokens to final outputs, showing every intermediate computation with concrete values.

## 9.1 Model Configuration and Initial Setup

Transformer specifications:

$$\text{vocab\_size} = 1000 \quad (73)$$

$$\text{d\_model} = 8 \text{ (embedding dimension)} \quad (74)$$

$$\text{n\_heads} = 2 \text{ (attention heads)} \quad (75)$$

$$\text{n\_kv\_heads} = 2 \text{ (key-value heads)} \quad (76)$$

$$\text{head\_dim} = 4 \text{ (per-head dimension)} \quad (77)$$

$$\text{seq\_len\_max} = 6 \quad (78)$$

$$\text{batch\_size} = 2 \quad (79)$$

Layer components:

- **Embedding layer:** Maps tokens to `d_model` dimensional vectors
- **Attention layer:** Multi-head self-attention with `CacheView`
- **Feed-forward layer:** Simple MLP transformation
- **Output projection:** Maps back to vocabulary space

## 9.2 Pass 1: Initial Processing (Prefill Phase)

Input tokens:

```
1 # Two sequences with different lengths
2 input_tokens = [
3     [45, 123, 67],          # Sequence 0: "Hello world example"
4     [89, 234]              # Sequence 1: "AI system"
5 ]
6 # Flattened for processing: [45, 123, 67, 89, 234]
7 # seqlens = [3, 2]
```

Step 1: Token Embedding

```
1 # Embedding lookup (simplified with concrete values)
2 embeddings = [
3     # Sequence 0 embeddings
4     [1.0, 0.5, -0.2, 0.8, 1.2, -0.5, 0.3, 0.9],    # token 45: "Hello"
5     [0.2, 1.1, 0.7, -0.3, 0.6, 1.0, -0.8, 0.4],    # token 123: "world"
6     [-0.5, 0.8, 1.3, 0.2, -0.7, 0.5, 1.1, -0.2],   # token 67: "example"
7     # Sequence 1 embeddings
8     [0.9, -0.4, 0.6, 1.1, 0.3, -0.9, 0.7, 0.8],    # token 89: "AI"
9     [1.1, 0.3, -0.6, 0.4, 0.9, 0.2, -0.4, 1.2]      # token 234: "system"
10 ]
11 # Shape: (5, 8) - 5 total tokens, 8-dimensional embeddings
```

Step 2: Positional Encoding

```

1 # Add positional encodings based on absolute positions
2 pos_encodings = [
3     [0.0, 1.0, 0.0, 1.0, 0.0, 1.0, 0.0, 1.0],      # Position 0
4     [0.1, 0.9, 0.1, 0.9, 0.1, 0.9, 0.1, 0.9],      # Position 1
5     [0.2, 0.8, 0.2, 0.8, 0.2, 0.8, 0.2, 0.8],      # Position 2
6     [0.0, 1.0, 0.0, 1.0, 0.0, 1.0, 0.0, 1.0],      # Position 0 (seq 1)
7     [0.1, 0.9, 0.1, 0.9, 0.1, 0.9, 0.1, 0.9]       # Position 1 (seq 1)
8 ]
9
10 # Add positional encodings to embeddings
11 input_embeddings = embeddings + pos_encodings
12 final_input = [
13     [1.0, 1.5, -0.2, 1.8, 1.2, 0.5, 0.3, 1.9],    # "Hello" + pos_0
14     [0.3, 2.0, 0.8, 0.6, 0.7, 1.9, -0.7, 1.3],    # "world" + pos_1
15     [-0.3, 1.6, 1.5, 1.0, -0.5, 1.3, 1.3, 0.6],   # "example" + pos_2
16     [0.9, 0.6, 0.6, 2.1, 0.3, 0.1, 0.7, 1.8],    # "AI" + pos_0
17     [1.2, 1.2, -0.5, 1.3, 1.0, 1.1, -0.3, 2.1]    # "system" + pos_1
18 ]

```

### Step 3: Multi-Head Attention - Query, Key, Value Projections

```

1 # Weight matrices (simplified 8x8 matrices, shown as key elements)
2 W_q = [
3     [0.1, 0.2, ..., 0.8],  # 8x8 query projection matrix
4     ...
5 ]
6 W_k = [
7     [0.2, 0.1, ..., 0.7],  # 8x8 key projection matrix
8     ...
9 ]
10 W_v = [
11     [0.3, 0.4, ..., 0.6],  # 8x8 value projection matrix
12     ...
13 ]
14
15 # Project to Q, K, V (showing calculated results)
16 Q = [
17     # Head 0 (dim 0-3), Head 1 (dim 4-7)
18     [2.1, 1.5, 0.8, 1.2, 1.8, 0.9, 1.1, 1.4],    # Query for "Hello"
19     [1.8, 2.2, 1.1, 0.9, 1.5, 1.7, 0.8, 1.6],    # Query for "world"
20     [1.4, 1.9, 1.3, 1.0, 1.2, 1.4, 1.0, 1.1],    # Query for "example"
21     [2.0, 1.3, 0.9, 1.5, 1.6, 1.0, 1.3, 1.7],    # Query for "AI"
22     [1.7, 1.6, 1.2, 1.1, 1.4, 1.3, 1.1, 1.5]     # Query for "system"
23 ]
24
25 K = [
26     [1.9, 1.4, 0.7, 1.1, 1.6, 0.8, 1.0, 1.3],    # Key for "Hello"
27     [1.6, 2.0, 1.0, 0.8, 1.3, 1.5, 0.7, 1.4],    # Key for "world"
28     [1.2, 1.7, 1.1, 0.9, 1.0, 1.2, 0.9, 1.0],    # Key for "example"
29     [1.8, 1.2, 0.8, 1.3, 1.4, 0.9, 1.1, 1.5],    # Key for "AI"
30     [1.5, 1.4, 1.0, 1.0, 1.2, 1.1, 1.0, 1.3]     # Key for "system"
31 ]
32
33 V = [
34     [2.2, 1.6, 0.9, 1.3, 1.9, 1.0, 1.2, 1.5],    # Value for "Hello"

```

```

35 [1.9, 2.3, 1.2, 1.0, 1.6, 1.8, 0.9, 1.7],      # Value for "world"
36 [1.5, 2.0, 1.4, 1.1, 1.3, 1.5, 1.1, 1.2],      # Value for "example"
37 [2.1, 1.4, 1.0, 1.6, 1.7, 1.1, 1.4, 1.8],      # Value for "AI"
38 [1.8, 1.7, 1.3, 1.2, 1.5, 1.4, 1.2, 1.6]       # Value for "system"
39 ]

```

#### Step 4: Initialize CacheView (First Pass)

```

1 # Initialize cache for this layer
2 cache_k = torch.zeros(2, 6, 2, 4)    # (batch, max_seq, n_kv_heads, head_dim
3
4 cache_v = torch.zeros(2, 6, 2, 4)
5 kv_seqlens = torch.tensor([0, 0])    # Initially empty
6
7 # CacheView setup for prefill
8 metadata = CacheInputMetadata(
9     to_cache_mask = [True, True, True, True, True],   # Cache all tokens
10    cache_positions = [0, 1, 2, 0, 1],                 # Cache positions
11    seqlens = [3, 2],                                  # Sequence lengths
12    prefill = True
13 )
14 cache_view = CacheView(cache_k, cache_v, metadata, kv_seqlens)

```

#### Step 5: Cache Update

```

1 # Reshape K, V for caching (split by heads)
2 K_reshaped = [
3     # Sequence 0
4     [[1.9, 1.4, 0.7, 1.1], [1.6, 0.8, 1.0, 1.3]],  # "Hello" - head 0,
5     # head 1
6     [[1.6, 2.0, 1.0, 0.8], [1.3, 1.5, 0.7, 1.4]],  # "world"
7     [[1.2, 1.7, 1.1, 0.9], [1.0, 1.2, 0.9, 1.0]],  # "example"
8     # Sequence 1
9     [[1.8, 1.2, 0.8, 1.3], [1.4, 0.9, 1.1, 1.5]],  # "AI"
10    [[1.5, 1.4, 1.0, 1.0], [1.2, 1.1, 1.0, 1.3]]   # "system"
11 ]
12
13 # Update cache using CacheView
14 cache_view.update(K_reshaped, V_reshaped)    # Similar structure for V
15
16 # Cache state after update:
17 # kv_seqlens = [3, 2]
18 # cache_k[0, 0:3] contains keys for sequence 0
19 # cache_k[1, 0:2] contains keys for sequence 1

```

#### Step 6: Attention Computation

```

1 # Compute attention scores (Q @ K^T)
2 # For sequence 0 (3x3 attention matrix for each head):
3 attention_scores_seq0_head0 = [
4     # Q_hello @ [K_hello, K_world, K_example]^T
5     [2.1*1.9 + 1.5*1.4 + 0.8*0.7 + 1.2*1.1,  # great->Hello
6      2.1*1.6 + 1.5*2.0 + 0.8*1.0 + 1.2*0.8,  # great->world
7      2.1*1.2 + 1.5*1.7 + 0.8*1.1 + 1.2*0.9], # great->example
8
9     # Q_world @ [K_hello, K_world, K_example]^T

```

```

10 [1.8*1.9 + 2.2*1.4 + 1.1*0.7 + 0.9*1.1, # world->Hello
11   1.8*1.6 + 2.2*2.0 + 1.1*1.0 + 0.9*0.8, # world->world
12   1.8*1.2 + 2.2*1.7 + 1.1*1.1 + 0.9*0.9], # world->example
13
14 # Q_example @ [K_hello, K_world, K_example]^T
15 [1.4*1.9 + 1.9*1.4 + 1.3*0.7 + 1.0*1.1, # example->Hello
16   1.4*1.6 + 1.9*2.0 + 1.3*1.0 + 1.0*0.8, # example->world
17   1.4*1.2 + 1.9*1.7 + 1.3*1.1 + 1.0*0.9] # example->example
18 ]
19 # = [9.15, 8.80, 8.45, 8.93]
20
21 # Apply causal mask (lower triangular)
22 masked_scores_seq0_head0 = [9.15, 8.80, 8.45, 8.93] # Can attend to all
23 previous
24 softmax_scores_seq0_head0 = [0.28, 0.24, 0.21, 0.27] # Normalized
25 attention
26
27 # Weighted combination of all cached values
28 great_output = 0.28*V_hello + 0.24*V_world + 0.21*V_example + 0.27*V_great

```

### Step 7: Apply Attention to Values

```

1 # Weighted sum of values using attention scores
2 attention_output_seq0_head0 = [
3     # Hello output (only attends to itself)
4     1.0 * [2.2, 1.6, 0.9, 1.3] = [2.2, 1.6, 0.9, 1.3],
5
6     # World output (0.27 * Hello + 0.73 * world)
7     0.27 * [2.2, 1.6, 0.9, 1.3] + 0.73 * [1.9, 2.3, 1.2, 1.0]
8     = [0.594, 0.432, 0.243, 0.351] + [1.387, 1.679, 0.876, 0.730]
9     = [1.981, 2.111, 1.119, 1.081],
10
11    # Example output (0.15*Hello + 0.52*world + 0.33*example)
12    0.15 * [2.2, 1.6, 0.9, 1.3] + 0.52 * [1.9, 2.3, 1.2, 1.0] + 0.33 *
13    [1.5, 2.0, 1.4, 1.1]
14    = [0.33, 0.24, 0.135, 0.195] + [0.988, 1.196, 0.624, 0.520] + [0.495,
15    0.66, 0.462, 0.363]
16    = [1.813, 2.096, 1.221, 1.078]
17 ]
18
19 # Similar computation for head 1 and sequence 1...

```

### Step 8: Final Output Projection

```

1 # Concatenate heads and project back to d_model
2 attention_output_concatenated = [
3     # Sequence 0
4     [2.2, 1.6, 0.9, 1.3, 1.8, 0.9, 1.1, 1.4], # Hello (head0 + head1)
5     [1.981, 2.111, 1.119, 1.081, 1.6, 1.8, 0.9, 1.7], # world
6     [1.813, 2.096, 1.221, 1.078, 1.4, 1.6, 1.2, 1.3], # example
7     # Sequence 1
8     [2.1, 1.4, 1.0, 1.6, 1.7, 1.1, 1.4, 1.8], # AI
9     [1.9, 1.7, 1.2, 1.4, 1.5, 1.3, 1.1, 1.5] # system
10 ]
11
12 # Feed-forward network processing (simplified)

```

```

13 ff_output = apply_feedforward(attention_output_concatenated)
14 pass1_output = [
15     [2.5, 1.8, 1.1, 1.6, 2.0, 1.2, 1.4, 1.7],      # Processed Hello
16     [2.3, 2.4, 1.3, 1.3, 1.8, 2.0, 1.1, 1.9],      # Processed world
17     [2.1, 2.3, 1.4, 1.3, 1.6, 1.8, 1.4, 1.5],      # Processed example
18     [2.4, 1.6, 1.2, 1.8, 1.9, 1.3, 1.6, 2.0],      # Processed AI
19     [2.2, 1.9, 1.4, 1.6, 1.7, 1.5, 1.3, 1.7]       # Processed system
20 ]

```

### 9.3 Pass 2: Incremental Generation

New tokens to generate:

```

1 # Generate next tokens for each sequence
2 new_tokens = [456, 789]    # "great", "works"
3 seqlens = [1, 1]           # One new token per sequence
4 # Current kv_seqlens = [3, 2] from previous pass

```

#### Step 1: Process New Token Embeddings

```

1 # New embeddings + positional encoding
2 new_embeddings = [
3     [1.3, 0.7, -0.4, 0.9, 1.1, -0.3, 0.5, 1.0],   # "great" + pos_3
4     [0.8, 1.2, 0.6, -0.2, 0.7, 1.1, -0.6, 0.9]     # "works" + pos_2
5 ]
6
7 # Project to Q, K, V
8 new_Q = [
9     [2.3, 1.7, 1.0, 1.4, 1.9, 1.1, 1.3, 1.6],      # Query for "great"
10    [1.9, 1.5, 1.1, 1.2, 1.6, 1.3, 1.0, 1.4]        # Query for "works"
11 ]
12
13 new_K = [
14     [2.0, 1.5, 0.9, 1.2, 1.7, 1.0, 1.1, 1.4],      # Key for "great"
15     [1.7, 1.3, 1.0, 1.1, 1.4, 1.2, 0.9, 1.3]        # Key for "works"
16 ]
17
18 new_V = [
19     [2.4, 1.8, 1.1, 1.5, 2.0, 1.2, 1.4, 1.7],      # Value for "great"
20     [2.0, 1.6, 1.3, 1.3, 1.7, 1.4, 1.1, 1.5]        # Value for "works"
21 ]

```

#### Step 2: Update Cache with New Tokens

```

1 # Update metadata for incremental generation
2 metadata_gen = CacheInputMetadata(
3     to_cache_mask = [True, True],          # Cache both new tokens
4     cache_positions = [3, 2],            # Next positions in cache
5     seqlens = [1, 1],                  # One token each
6     prefill = False                   # Generation phase
7 )
8
9 # Update cache using CacheView
10 new_K_reshaped = [
11     [[2.0, 1.5, 0.9, 1.2], [1.7, 1.0, 1.1, 1.4]], # "great"

```

```

12     [[1.7, 1.3, 1.0, 1.1], [1.4, 1.2, 0.9, 1.3]]    # "works"
13 ]
14
15 cache_view.update(new_K_reshaped, new_V_reshaped)
16 # kv_seqlens now = [4, 3]

```

### Step 3: Attention with Cached Context

```

1 # Use interleave_kv to combine cached and new keys/values
2 K_full, V_full = cache_view.interleave_kv(new_K_reshaped, new_V_reshaped)
3
4 # K_full now contains:
5 # Sequence 0: [K_hello, K_world, K_example, K_great] (4 tokens)
6 # Sequence 1: [K_AI, K_system, K_works] (3 tokens)
7
8 # Compute attention for new queries against full context
9 attention_scores_great = [
10     # Q_great @ [K_hello, K_world, K_example, K_great]^T
11     [2.3*1.9 + 1.7*1.4 + 1.0*0.7 + 1.4*1.1, # great->Hello
12      2.3*1.6 + 1.7*2.0 + 1.0*1.0 + 1.4*0.8, # great->world
13      2.3*1.2 + 1.7*1.7 + 1.0*1.1 + 1.4*0.9, # great->example
14      2.3*2.0 + 1.7*1.5 + 1.0*0.9 + 1.4*1.2] # great->great
15 ]
16 # = [9.15, 8.80, 8.45, 8.93]
17
18 # Apply causal mask and softmax
19 masked_scores_great = [9.15, 8.80, 8.45, 8.93] # Can attend to all
20 previous
21 softmax_great = [0.28, 0.24, 0.21, 0.27]           # Normalized attention
22
23 # Weighted combination of all cached values
24 great_output = 0.28*V_hello + 0.24*V_world + 0.21*V_example + 0.27*V_great

```

## 9.4 Pass 3: Continued Generation

Generate one more token:

```

1 # Generate final tokens
2 final_tokens = [512] # "!" for sequence 0 only
3 seqlens = [1, 0]      # Only sequence 0 gets new token
4
5 # Current kv_seqlens = [4, 3]
6 # After this pass: kv_seqlens = [5, 3]

```

Complete processing similar to Pass 2, showing cache growth:

```

1 # Final cache state:
2 # Sequence 0: [Hello, world, example, great, !] (5 tokens cached)
3 # Sequence 1: [AI, system, works] (3 tokens cached)
4
5 # Cache utilization:
6 # cache_k shape: (2, 6, 2, 4)
7 # Used slots: 5 + 3 = 8 out of 12 total slots = 67% utilization
8
9 # Total tokens processed: 3 + 2 + 1 + 1 + 1 = 8 tokens

```

```

10 # Cached tokens available: 5 + 3 = 8 tokens
11 # New computations needed: Only for new tokens (3 in total across passes)
12 # Computation savings: 8 cached / 3 new = 267% efficiency gain

```

## 9.5 Performance Summary Across 3 Passes

Memory efficiency evolution:

$$\text{Pass 1 (prefill): Cache 5 tokens, compute 5 tokens} \quad (80)$$

$$\text{Pass 2 (generation): Use 5 cached + cache 2 new = 7 total} \quad (81)$$

$$\text{Pass 3 (generation): Use 7 cached + cache 1 new = 8 total} \quad (82)$$

$$(83)$$

$$\text{Total computations without cache: } 5 + 7 + 8 = 20 \text{ token-computations} \quad (84)$$

$$\text{Total computations with cache: } 5 + 2 + 1 = 8 \text{ token-computations} \quad (85)$$

$$\text{Efficiency gain: } \frac{20}{8} = 2.5 \times \text{ speedup} \quad (86)$$

This comprehensive 3-pass example demonstrates how `CacheView` seamlessly integrates into the complete transformer pipeline, enabling efficient incremental generation while maintaining perfect attention semantics across all computation phases.

## 10 Visual Understanding: The Core Insight Behind CacheView

The fundamental insight that makes `CacheView` possible is beautifully illustrated in the inference optimization diagram. Let's break down this concept with concrete examples that directly relate to our cache implementation.

### 10.1 The Key Insight: Reusing Computed Attention Components

**The Problem:** During autoregressive generation, we repeatedly compute attention for tokens we've already processed.

**The Solution:** Cache the Key (K) and Value (V) matrices since they don't change for previously processed tokens.

### 10.2 Detailed Walkthrough: From Theory to CacheView Implementation

Let's trace through the exact scenario shown in the diagram using our `CacheView` implementation:

**Setup:** Generating token 4 in a sequence that already has tokens 1, 2, and 3.

```

1 # Initial state: we have processed 3 tokens
2 existing_tokens = ["The", "cat", "sat"] # tokens 1, 2, 3
3 current_seqlens = [3] # sequence length so far
4 vocab_size = 50000
5 d_model = 4096
6 n_heads = 32
7 head_dim = 128 # d_model / n_heads

```

### Step 1: The Inefficient Approach (Without Caching)

```

1 # Without caching, we would recompute everything for token 4
2 all_tokens = ["The", "cat", "sat", "on"] # including new token 4
3 seqlen = 4
4
5 # Compute Q, K, V for ALL tokens (wasteful!)
6 Q_all = compute_queries(all_tokens) # Shape: (4, 4096)
7 K_all = compute_keys(all_tokens) # Shape: (4, 4096)
8 V_all = compute_values(all_tokens) # Shape: (4, 4096)
9
10 # Attention computation
11 attention_scores = Q_all @ K_all.T # Shape: (4, 4)
12 # Only the last row (token 4's attention) is actually needed!

```

### Step 2: The CacheView Approach (Efficient)

```

1 # Initialize CacheView with pre-allocated cache
2 cache_k = torch.zeros(1, 8, 32, 128) # (batch=1, max_seq=8, heads=32,
3   head_dim=128)
4 cache_v = torch.zeros(1, 8, 32, 128)
5 kv_seqlens = torch.tensor([3]) # We have 3 tokens cached
6
7 # Previous tokens' K, V are already cached from earlier computations
8 # cached_k contains: [K_The, K_cat, K_sat] at positions [0, 1, 2]
9 # cached_v contains: [V_The, V_cat, V_sat] at positions [0, 1, 2]
10
11 # For token 4 "on", we only compute NEW Q, K, V
12 new_token = "on"
13 Q_new = compute_queries([new_token]) # Shape: (1, 4096) - only 1 token
14   !
15 K_new = compute_keys([new_token]) # Shape: (1, 4096) - only 1 token
16   !
17 V_new = compute_values([new_token]) # Shape: (1, 4096) - only 1 token
18   !
19
20 # Update cache with new K, V
21 metadata = CacheInputMetadata(
22     to_cache_mask=[True], # Cache the new token
23     cache_positions=[3], # Store at position 3
24     seqlens=[1], # 1 new token
25     prefill=False # Generation phase
26 )
27
28 cache_view = CacheView(cache_k, cache_v, metadata, kv_seqlens)
29 K_new_reshaped = K_new.view(1, 32, 128) # Reshape for heads
30 V_new_reshaped = V_new.view(1, 32, 128)
31 cache_view.update(K_new_reshaped, V_new_reshaped)

```

### Step 3: Efficient Attention Computation

```
1 # Get complete K, V using interleave_kv
2 K_complete, V_complete = cache_view.interleave_kv(K_new_reshaped,
3     V_new_reshaped)
4 # K_complete now contains: [K_The, K_cat, K_sat, K_on]
5 # V_complete now contains: [V_The, V_cat, V_sat, V_on]
6
7 # Compute attention ONLY for the new token's query
8 Q_new_reshaped = Q_new.view(1, 32, 128)
9 attention_scores = Q_new_reshaped @ K_complete.transpose(-2, -1) # Shape:
10    (1, 32, 4)
11
12 # Apply causal mask (token 4 can attend to tokens 1, 2, 3, 4)
13 causal_mask = torch.tril(torch.ones(1, 4)) # Lower triangular
14 masked_scores = attention_scores.masked_fill(causal_mask == 0, float('-inf'))
15
16 # Softmax and weighted sum
17 attention_weights = torch.softmax(masked_scores / math.sqrt(128), dim=-1)
18 output = attention_weights @ V_complete # Final attention output for
19      token 4
```

### 10.3 Mapping to the Visual Diagram

Let's directly map our implementation to each numbered insight in the diagram:

**Insight 1: "We already computed these dot products. Can we cache them?"**

```
1 # In our CacheView implementation:
2 # - Previous Q @ K^T computations are implicitly "cached" by caching K
3 # - When we need Q_new @ K_old^T, we retrieve K_old from cache
4 # - No need to recompute K for tokens 1, 2, 3
5
6 # Example: Token 4 attending to token 2
7 # Instead of recomputing: Q_on @ K_cat^T
8 # We use: Q_on @ cached_K[1]^T (where cached_K[1] is K_cat)
```

**Insight 2: "Since the model is causal, we don't care about attention of a token with its successors"**

```
1 # CacheView respects causality automatically
2 # When generating token 4, the attention computation is:
3
4 attention_pattern = [
5     # Token 4 can attend to:
6     [True, True, True, True],    # tokens 1, 2, 3, 4 (all previous +
7     self)
8     # Tokens 1, 2, 3 don't need to attend to token 4 (future)
9     # We don't compute their attention - they're already finalized!
10 ]
11
12 # This is why we only compute Q for the NEW token
13 # Previous tokens' outputs are already computed and cached
```

**Insight 3: "We don't care about these, as we want to predict the next token"**

```

1 # In generation, we only need the output for the LAST token
2 # Previous tokens' outputs are already computed
3
4 # CacheView optimization:
5 # - Don't recompute attention for tokens 1, 2, 3
6 # - Only compute attention for token 4
7 # - Previous hidden states are already available from earlier forward
    passes

```

**Insight 4: "We are only interested in this last row!"**

```

1 # This is exactly what CacheView enables:
2 # Instead of computing full (4 x 4) attention matrix:
3 full_attention = Q_all @ K_all.T # Shape: (4, 4) - wasteful!
4
5 # We compute only the last row:
6 last_row_attention = Q_new @ K_complete.T # Shape: (1, 4) - efficient!
7
8 # The last row is: [Q_"on" @ K_"The"^-T, Q_"on" @ K_"cat"^-T, Q_"on" @ K_"
    sat"^-T, Q_"on" @ K_"on"^-T]
9 # Where K_"The", K_"cat", K_"sat" come from cache, and K_"on" is newly
    computed

```

## 10.4 Concrete Performance Comparison

Without CacheView (Naive Approach):

Token 4 generation: (87)

Q computation:  $4 \text{ tokens} \times 4096 \times 4096 = 67M$  operations (88)

K computation:  $4 \text{ tokens} \times 4096 \times 4096 = 67M$  operations (89)

V computation:  $4 \text{ tokens} \times 4096 \times 4096 = 67M$  operations (90)

Attention:  $4 \times 4 \times 4096 = 65K$  operations (91)

Total:  $\approx 201M$  operations (92)

With CacheView (Optimized):

Token 4 generation: (93)

Q computation:  $1 \text{ token} \times 4096 \times 4096 = 17M$  operations (94)

K computation:  $1 \text{ token} \times 4096 \times 4096 = 17M$  operations (95)

V computation:  $1 \text{ token} \times 4096 \times 4096 = 17M$  operations (96)

Attention:  $1 \times 4 \times 4096 = 16K$  operations (97)

Cache operations:  $\approx 1K$  operations (98)

Total:  $\approx 51M$  operations (99)

**Speedup:**  $\frac{201M}{51M} = 3.9 \times$  faster!

## 10.5 Real-World Example: Chatbot Response Generation

Let's see how this applies to a practical chatbot scenario:

```

1 # User input: "What is machine learning?"
2 # Model generates: "Machine learning is a subset of artificial
3 # intelligence..."
4
5 conversation_so_far = [
6     "What", "is", "machine", "learning", "?" # User input (5 tokens)
7 ]
8
9 # Without CacheView: For each generated token, recompute everything
10 generated_tokens = []
11 for i in range(10): # Generate 10 response tokens
12     # Inefficient: recompute Q, K, V for ALL tokens so far
13     all_tokens = conversation_so_far + generated_tokens
14     # Cost grows quadratically: O(n2) where n is sequence length
15
16 # With CacheView: Incremental generation
17 cache_view = initialize_cache_view()
18
19 # Cache user input once
20 cache_user_input(conversation_so_far, cache_view)
21
22 # Generate response incrementally
23 for i in range(10):
24     # Efficient: only compute Q, K, V for NEW token
25     new_token = generate_next_token(cache_view)
26     cache_view.update(new_token) # Add to cache
27     # Cost is constant: O(1) per token
28 # Result: 10x faster response generation for long conversations!

```

## 10.6 The Mathematical Beauty

The core insight can be expressed mathematically:

**Traditional Approach:**

$$\text{Attention}_t = \text{softmax} \left( \frac{Q_{1:t} K_{1:t}^T}{\sqrt{d_k}} \right) V_{1:t}$$

Where we recompute  $Q_{1:t}$ ,  $K_{1:t}$ ,  $V_{1:t}$  for every new token  $t$ .

**CacheView Approach:**

$$\text{Attention}_t = \text{softmax} \left( \frac{[Q_t][\text{cached}(K_{1:t-1}), K_t]^T}{\sqrt{d_k}} \right) [\text{cached}(V_{1:t-1}), V_t]$$

Where: -  $Q_t$  is computed only for the new token -  $K_{1:t-1}$  and  $V_{1:t-1}$  are retrieved from cache - Only  $K_t$  and  $V_t$  are newly computed

This transforms the complexity from  $O(t^2)$  to  $O(t)$  for sequence generation!

## 10.7 Conclusion: From Insight to Implementation

The visual diagram captures the essential insight that makes modern transformer inference efficient. CacheView is the practical implementation of this insight, providing:

- **Automatic caching** of K and V matrices
- **Efficient retrieval** during attention computation
- **Seamless integration** with existing transformer code
- **Dramatic speedups** for autoregressive generation

The beauty lies in the simplicity: by recognizing that causal attention only needs the "last row" of the attention matrix, we can cache all previous computations and achieve massive efficiency gains without any loss in model quality.