

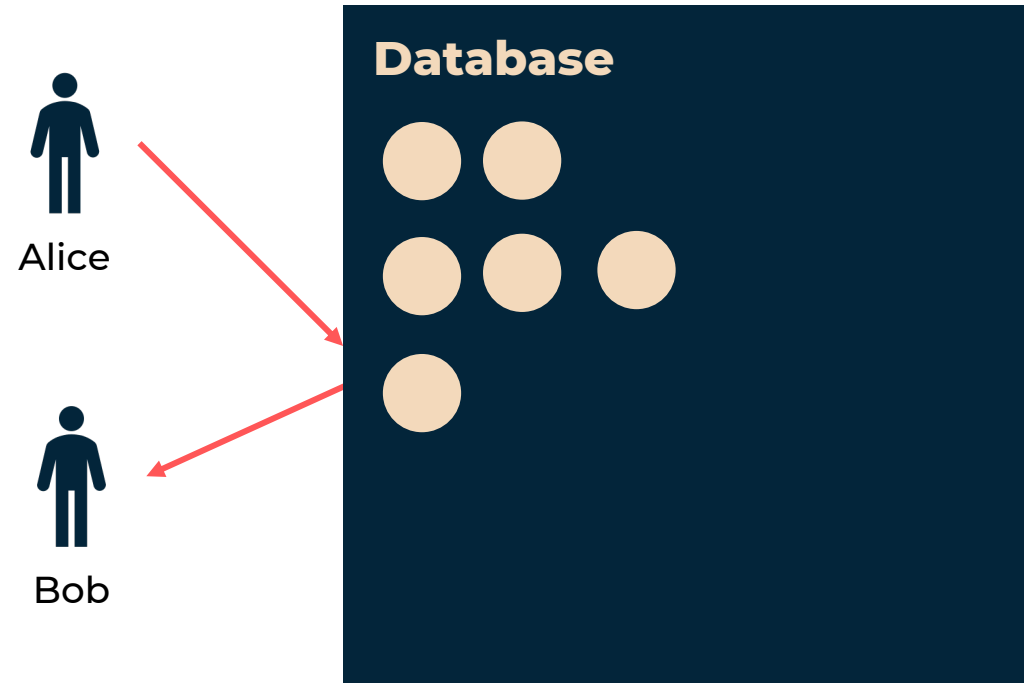
Towards correct high-performance database backend

Saalik Hatia - Annette Bieniusa - Carla Ferreira - Gustavo Petri - Marc Shapiro

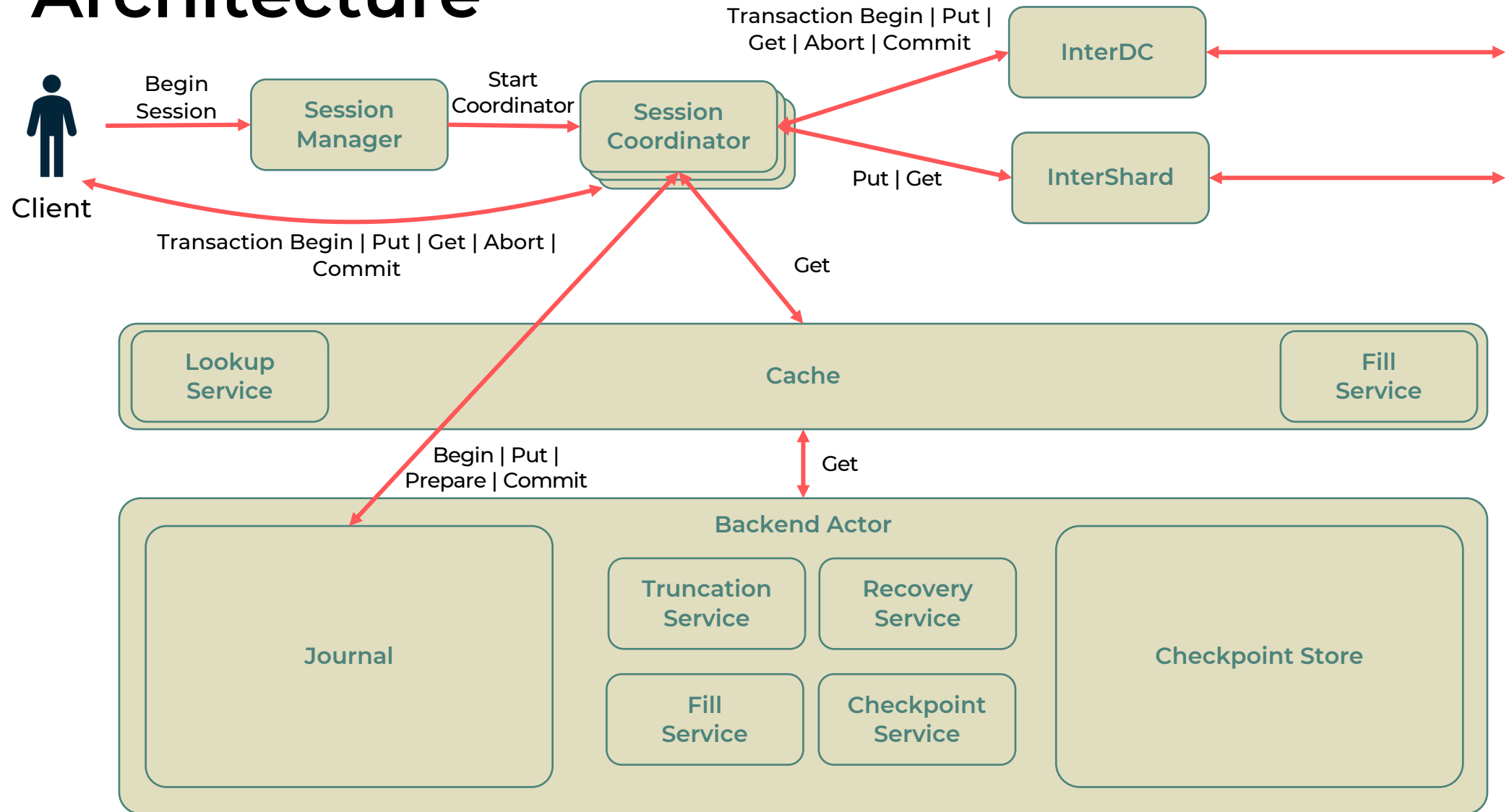


What is a database ?

- Store data
- Retrieve data
- Share data



Architecture



Target database - Antidote

Geo-distributed data centers (DC)

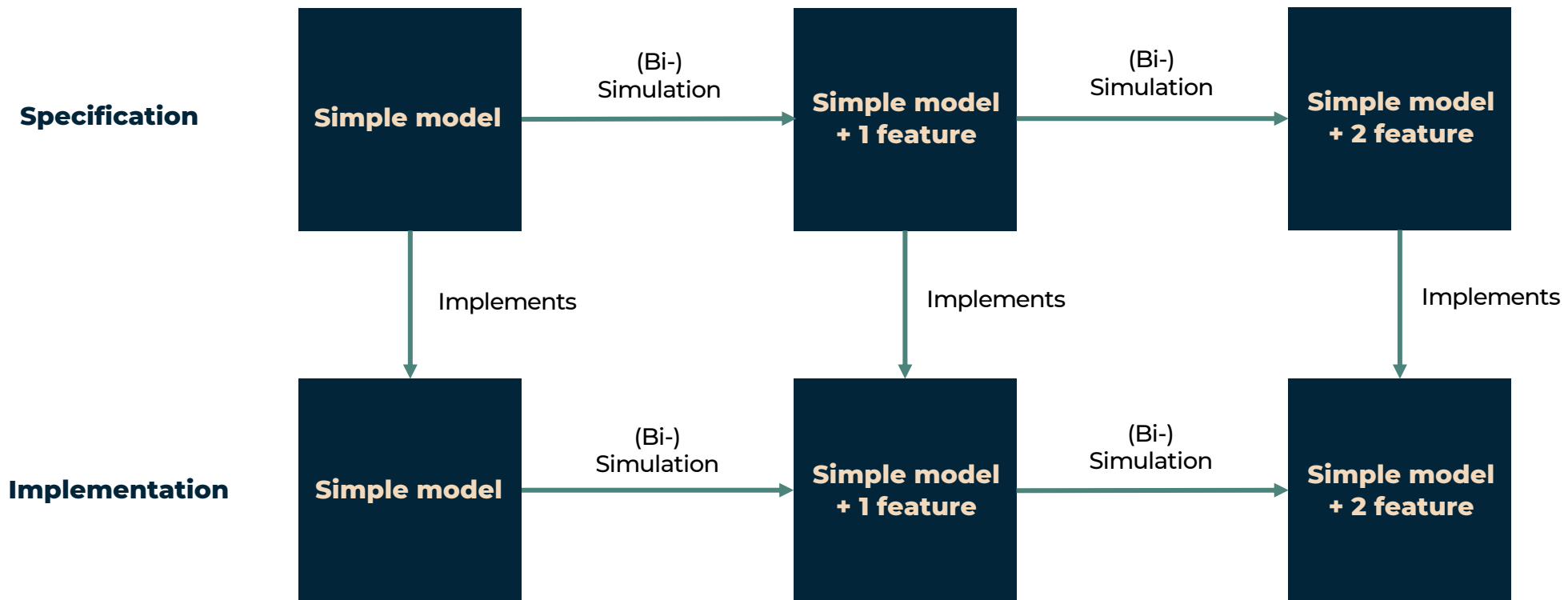
Full replication

Causal consistency between DCs

In a DC:

- Total order per DC [*Akkoorath et al. '16*]
- Sharding [*DeCandia et al '07*]
- Multiversion Concurrency Control (MVCC) with partial order
- Journal of operations
- Safe journal truncation
- States stored in a Checkpoint Store

Adding one feature at a time



Our plan

Baseline:

- Transaction
- Causal consistency
- Multi-version concurrency control
- Store all versions

Incrementally add features:

- Bounded versions
- Fault tolerance
- Journal
- Etc..

Specification of the model

Operational semantics

Invariants

$$\begin{array}{c}
 \text{READ LOCAL} \quad \frac{k \in \text{dom}(\mathcal{E}) \quad \mathcal{E}(k) = v}{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, \perp) \xrightarrow{\text{read}_\tau(k, \text{dt}) \Rightarrow \{v\}} (\tau, \text{dt}, \mathcal{E}, \mathcal{R}, \perp)} \quad \text{READ} \quad \frac{k \notin \text{dom}(\mathcal{E})}{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, \perp) \xrightarrow{\text{read}_\tau(k, \text{dt}) \Rightarrow V} (\tau, \text{dt}, \mathcal{E}, \mathcal{R} \cup \{k\}, \perp)} \\
 \text{EFFECT} \quad \frac{}{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, \perp) \xrightarrow{\text{effect}_\tau(k, v)} (\tau, \text{dt}, \mathcal{E}[k \leftarrow v], \mathcal{R}, \perp)}
 \end{array}$$

Figure 5: Transaction Semantics

$$\begin{array}{c}
 \text{STEP} \quad \frac{\mathcal{T} = \mathcal{T}'' \cup \{\tau\} \quad \tau \xrightarrow{\alpha_\tau} \tau' \quad \mathcal{T}' = \mathcal{T}'' \cup \{\tau'\}}{(\mathcal{T}_a, \mathcal{T}_c, \mathcal{T}) \xrightarrow[\tau]{\alpha} (\mathcal{T}_a, \mathcal{T}_c, \mathcal{T}')} \quad \text{ABORT} \quad \frac{\mathcal{T} = \mathcal{T}' \cup \{\tau\} \quad \mathcal{T}'_a = \mathcal{T}_a \cup \{\tau\}}{(\mathcal{T}_a, \mathcal{T}_c, \mathcal{T}) \xrightarrow[\tau]{\text{abort}_\tau} (\mathcal{T}'_a, \mathcal{T}_c, \mathcal{T}')} \\
 \text{BEGIN TRANSACTION} \quad \frac{\tau \notin \mathcal{T}_a \cup \mathcal{T}_c \cup \mathcal{T} \quad \mathcal{CM}_{\text{begin}}(\text{dt}, \neg, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}) \quad \mathcal{T}' = \mathcal{T} \cup \{(\tau, \text{dt}, \emptyset, \emptyset, \perp)\}}{(\mathcal{T}_a, \mathcal{T}_c, \mathcal{T}) \xrightarrow[\tau]{\text{begin}(\text{dt})} (\mathcal{T}_a, \mathcal{T}_c, \mathcal{T}')} \\
 \text{COMMIT} \quad \frac{\mathcal{T} = \mathcal{T}' \cup \{\tau\} \quad \mathcal{CM}_{\text{commit}}(\tau, \neg, \mathcal{T}_c, \text{ct}) \quad \tau = (\neg, \neg, \mathcal{E}, \neg, \neg) \quad \mathcal{T}'_c = \mathcal{T}_c \cup \{\tau[\perp \leftarrow \text{ct}]\}}{(\mathcal{T}_a, \mathcal{T}_c, \mathcal{T}) \xrightarrow[\tau]{\text{commit}_\tau(\text{ct}, \mathcal{E})} (\mathcal{T}_a, \mathcal{T}'_c, \mathcal{T}')}
 \end{array}$$

Figure 6: Transaction System

Starting simple

Client API

- *begin(dt)* initialize a transaction with a snapshot
- *read(key)* blob returns the object key from the snapshot Dt
- *effect(key, value)* assigns the value blob to the key
- *commit()* assign a commit timestamp to the transaction
- *abort()* abort the live transaction

Implementation

We take our specification to be as close as possible to the model

We represent a transaction with the following pattern:

$$(\tau_i, \mathbf{dt}, \mathcal{E}, \mathcal{R}, \mathbf{ct})$$

with the following interpretation:

- τ_i : transaction identifier (unique by definition),
- \mathbf{dt} : dependency timestamp (identifies the transactions this one depends upon),
- \mathcal{E} : effect map (records the content of the transaction's writes),
- \mathcal{R} : read set (records what objects the transaction has read),
- \mathbf{ct} : commit timestamp (set at termination time). In the case where the transaction has not been finalized we have that $\mathbf{ct} = \perp$.

Implementation

We take our specification to be as close as possible to the model

```
public class Transaction {  
    // Unique by definition  
    private TransactionID id;  
    // Identifies the transactions this one depends upon  
    private Timestamps dependency;  
    // Records the content of the transaction's writes  
    private HashMap<String, Value> effectMap;  
    // Records what objects the transaction has read  
    private CopyOnWriteArraySet<String> readSet;  
    // Time of commit  
    private Timestamps commit;  
}
```

Implementation

Then we add all the necessary invariants

Begin transaction:

- Only one active transaction per client
- Dependency timestamp must be valid

Implementation

Then we add all the necessary invariants.

```
@Override
public void startTransaction(TransactionID dependency) {
    checkArgument(tr == null, "Transaction already started");
    if (backend.dependencyIsValid(dependency)){
        tr = new Transaction(backend, dependency);
    }else{
        System.out.println("Dependency is not valid ! " +
            "Please retry with a correct transaction identifier");
    }
}
```

Implementation

Then we add all the necessary invariants.

```
@Override
public void startTransaction(TransactionID dependency) {
    checkArgument(tr == null, "Transaction already started");
    if (backend.dependencyIsValid(dependency)){
        tr = new Transaction(backend, dependency);
    }else{
        System.out.println("Dependency is not valid ! " +
            "Please retry with a correct transaction identifier");
    }
}
```

Implementation

Then we add all the necessary invariants.

```
@Override
public void startTransaction(TransactionID dependency) {
    checkArgument(tr == null, "Transaction already started");
    if (backend.dependencyIsValid(dependency)){
        tr = new Transaction(backend, dependency);
    }else{
        System.out.println("Dependency is not valid ! " +
            "Please retry with a correct transaction identifier");
    }
}
```

Adding a feature – Bounded Versions

Once everything is working, we add a single feature:

- Bounded versions
 - Write the specification
 - Implement the new version
 - Comparison with the previous model

Bounded Versions – Memory invariants

We define an arbitrary limitation on the size of the memory used by the system.

New system values:

$$\begin{array}{c} M_{used} \\ M_{limit} \end{array}$$

New invariant:

$$M_{used} \leq M_{limit}$$

Bounded Versions – Memory invariants

```
public void effect(String key, Value value) {
    checkArgument(backend.sizeAvailable(), "Memory size limit reached");

    @Nullable
    Value newValue = null;
    try {
        if (effectMap.containsKey(key)) {
            Value oldValue = effectMap.get(key);
            newValue = Value.merge(oldValue, value);
            effectMap.put(key, newValue);
        } else {
            Bounded.Value oldValue = backend.getValue(key, dependency);
            if (oldValue != null) {
                newValue = Value.merge(oldValue, value);
                effectMap.put(key, newValue);
            } else {
                effectMap.put(key, new Value(id, value));
            }
        }
    } finally {
        checkArgument(effectMap.containsValue(newValue), "Value was not added to effectmap");
    }
}
```

Bounded Versions – Memory invariants

```
public void effect(String key, Value value) {
    checkArgument(backend.sizeAvailable(), "Memory size limit reached");

    @Nullable
    Value newValue = null;
    try {
        if (effectMap.containsKey(key)) {
            Value oldValue = effectMap.get(key);
            newValue = Value.merge(oldValue, value);
            effectMap.put(key, newValue);
        } else {
            Bounded.Value oldValue = backend.getValue(key, dependency);
            if (oldValue != null) {
                newValue = Value.merge(oldValue, value);
                effectMap.put(key, newValue);
            } else {
                effectMap.put(key, new Value(id, value));
            }
        }
    } finally {
        checkArgument(effectMap.containsValue(newValue), "Value was not added to effectmap");
    }
}
```

Bounded Versions – Memory invariants

```
public void effect(String key, Value value) {
    checkArgument(backend.sizeAvailable(), "Memory size limit reached");

    @Nullable
    Value newValue = null;
    try {
        if (effectMap.containsKey(key)) {
            Value oldValue = effectMap.get(key);
            newValue = Value.merge(oldValue, value);
            effectMap.put(key, newValue);
        } else {
            Bounded.Value oldValue = backend.getValue(key, dependency);
            if (oldValue != null) {
                newValue = Value.merge(oldValue, value);
                effectMap.put(key, newValue);
            } else {
                effectMap.put(key, new Value(id, value));
            }
        }
    } finally {
        checkArgument(effectMap.containsValue(newValue), "Value was not added to effectmap");
    }
}
```

Bounded Versions – Memory invariants

```
public void effect(String key, Value value) {
    checkArgument(backend.sizeAvailable(), "Memory size limit reached");

    @Nullable
    Value newValue = null;
    try {
        if (effectMap.containsKey(key)) {
            effectMap.put(key, newValue);
        } else {
            effectMap.put(key, new Value(id, value));
        }
    } finally {
        checkArgument(effectMap.containsValue(newValue), "Value was not added to effectmap");
    }
}
```

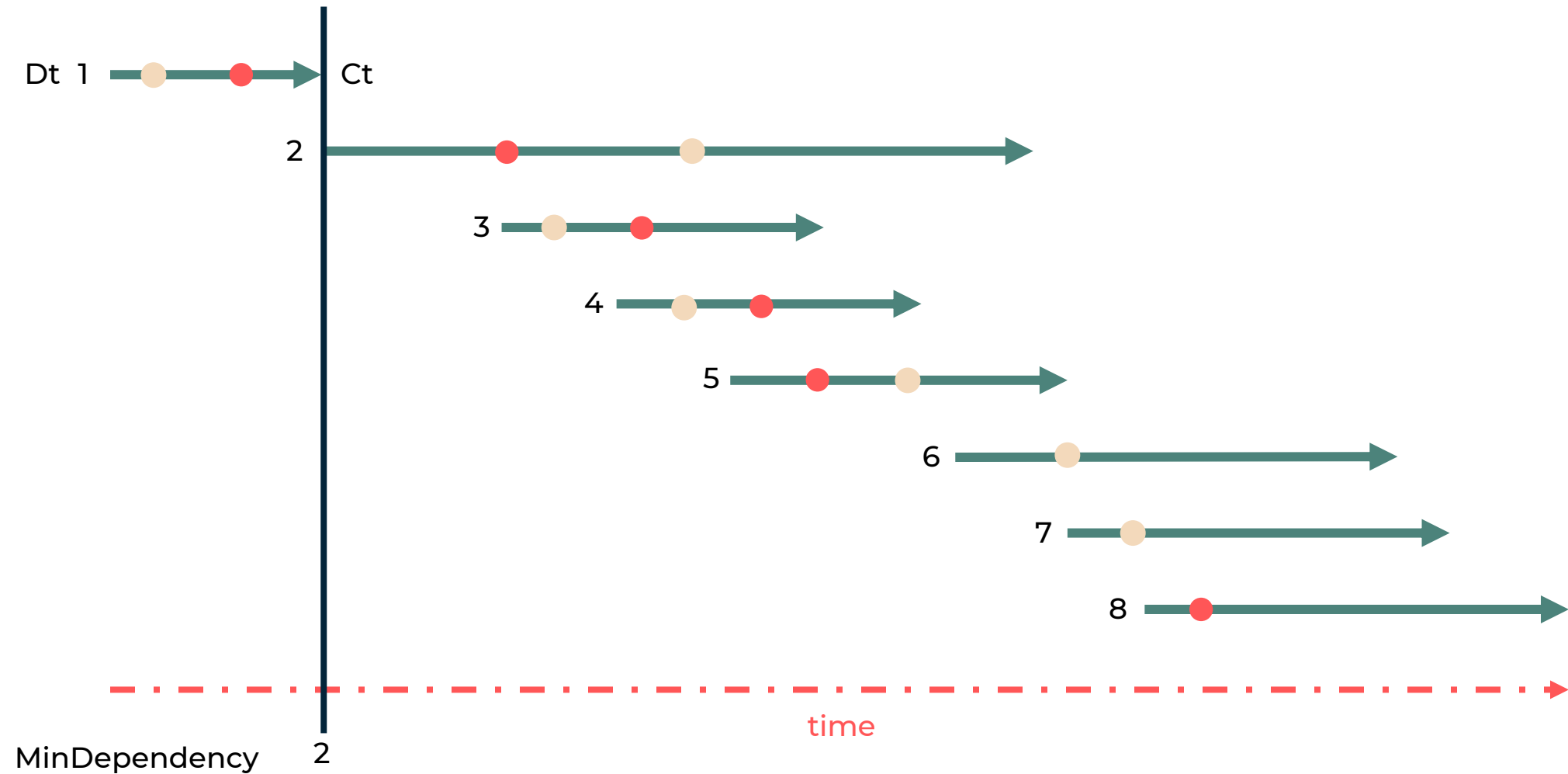
```
public boolean sizeAvailable() {
    int mUsed = backend.size();
    return mUsed <= MLIMIT;
}
```

Bounded memory – Garbage collection

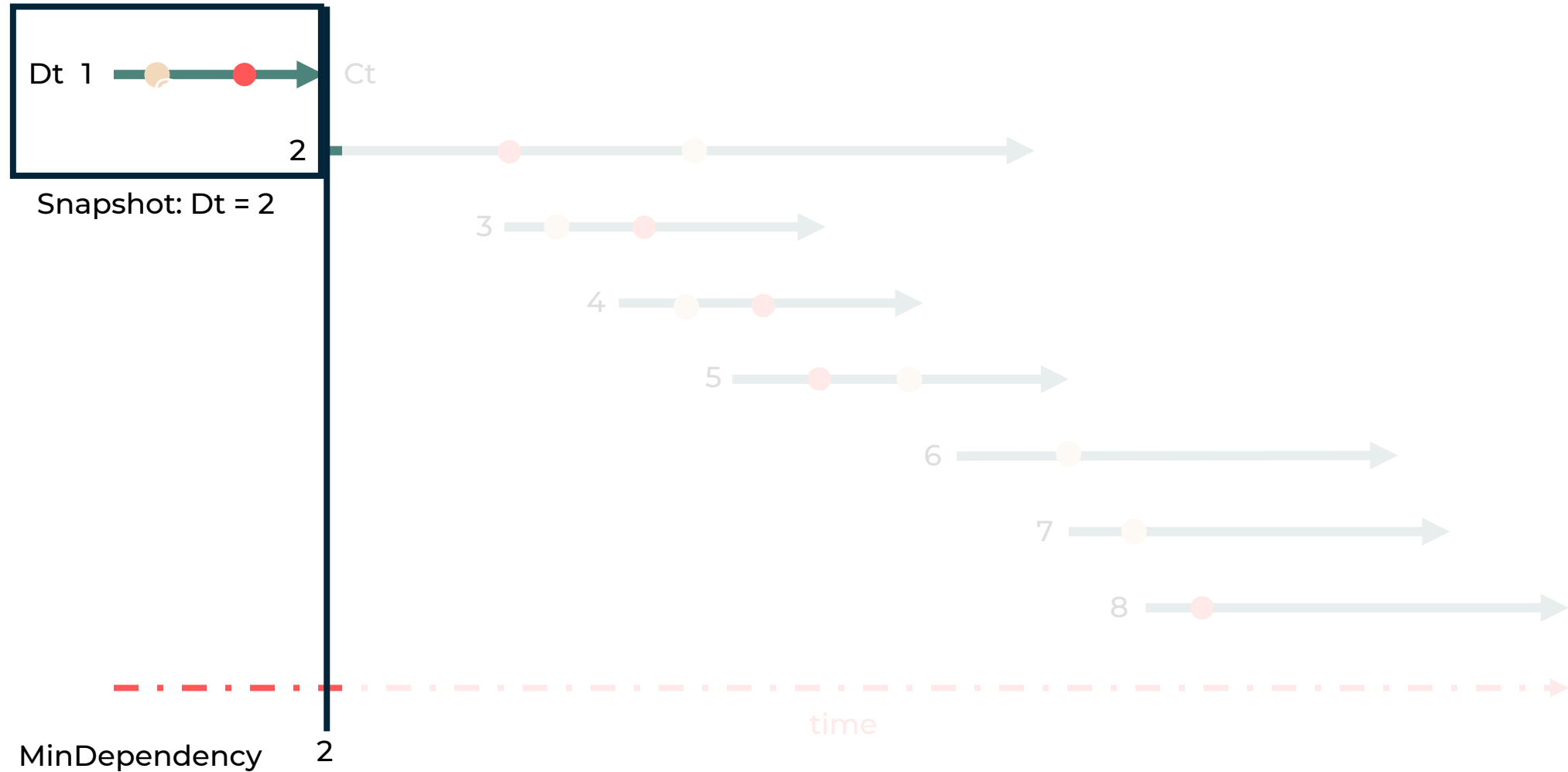
How do we know what we can throw away ?

- Keep track of the lowest dependency of a running transaction (MinDependency)
 - Among the versions contained in the snapshot
 - At the minimum, keep the most recent version of every key

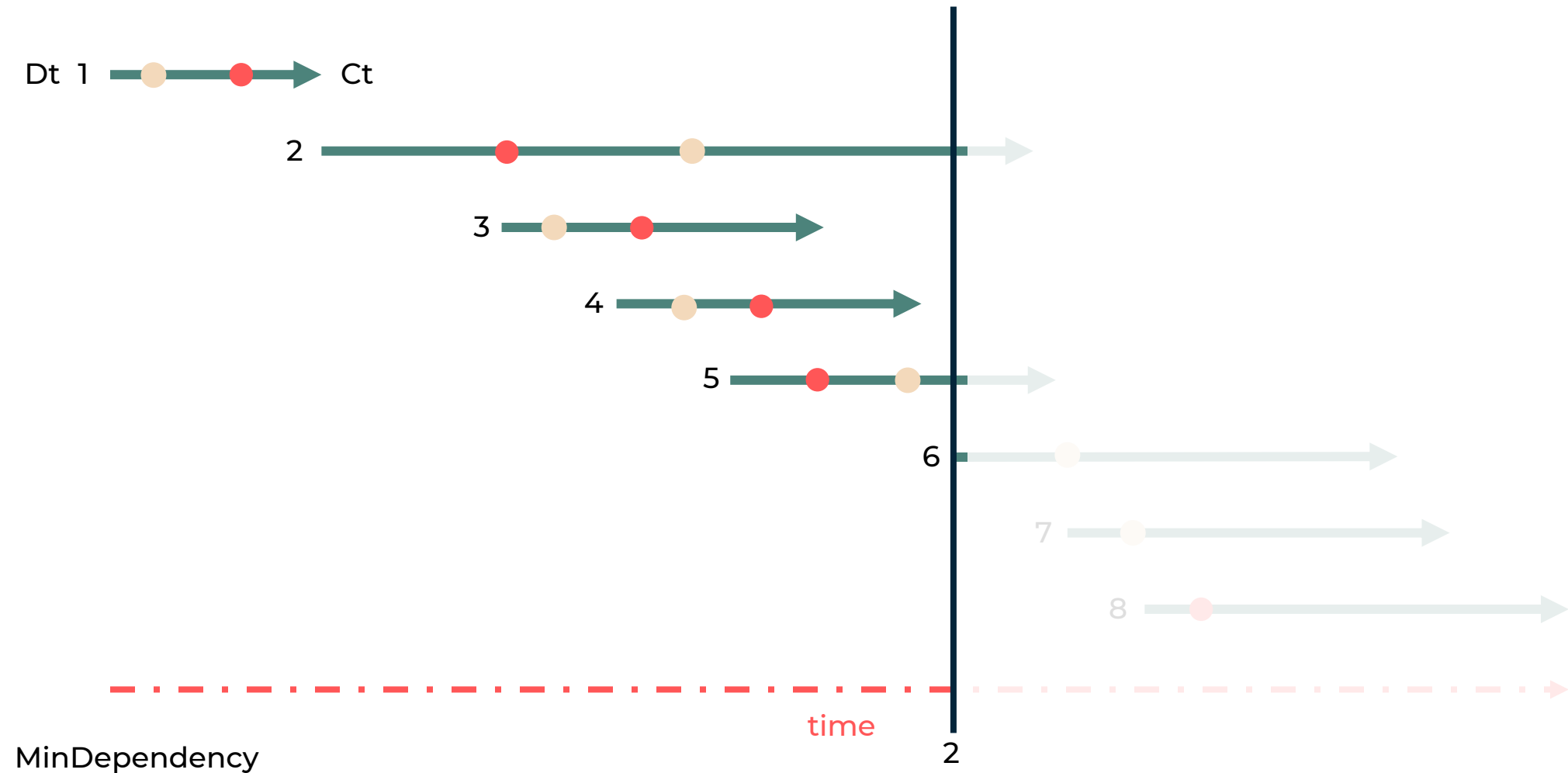
Bounded memory – Garbage collection



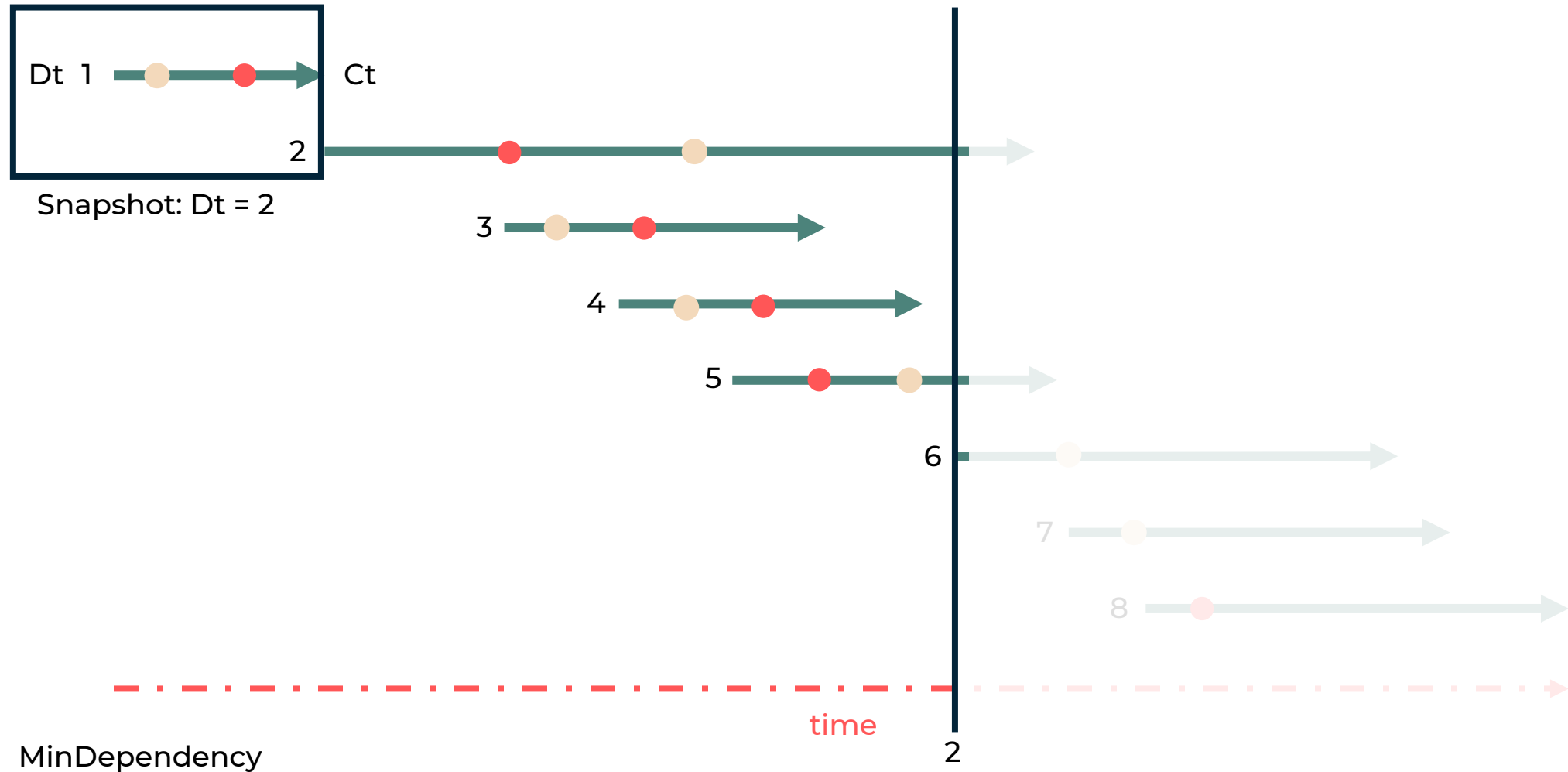
Bounded memory – Garbage collection



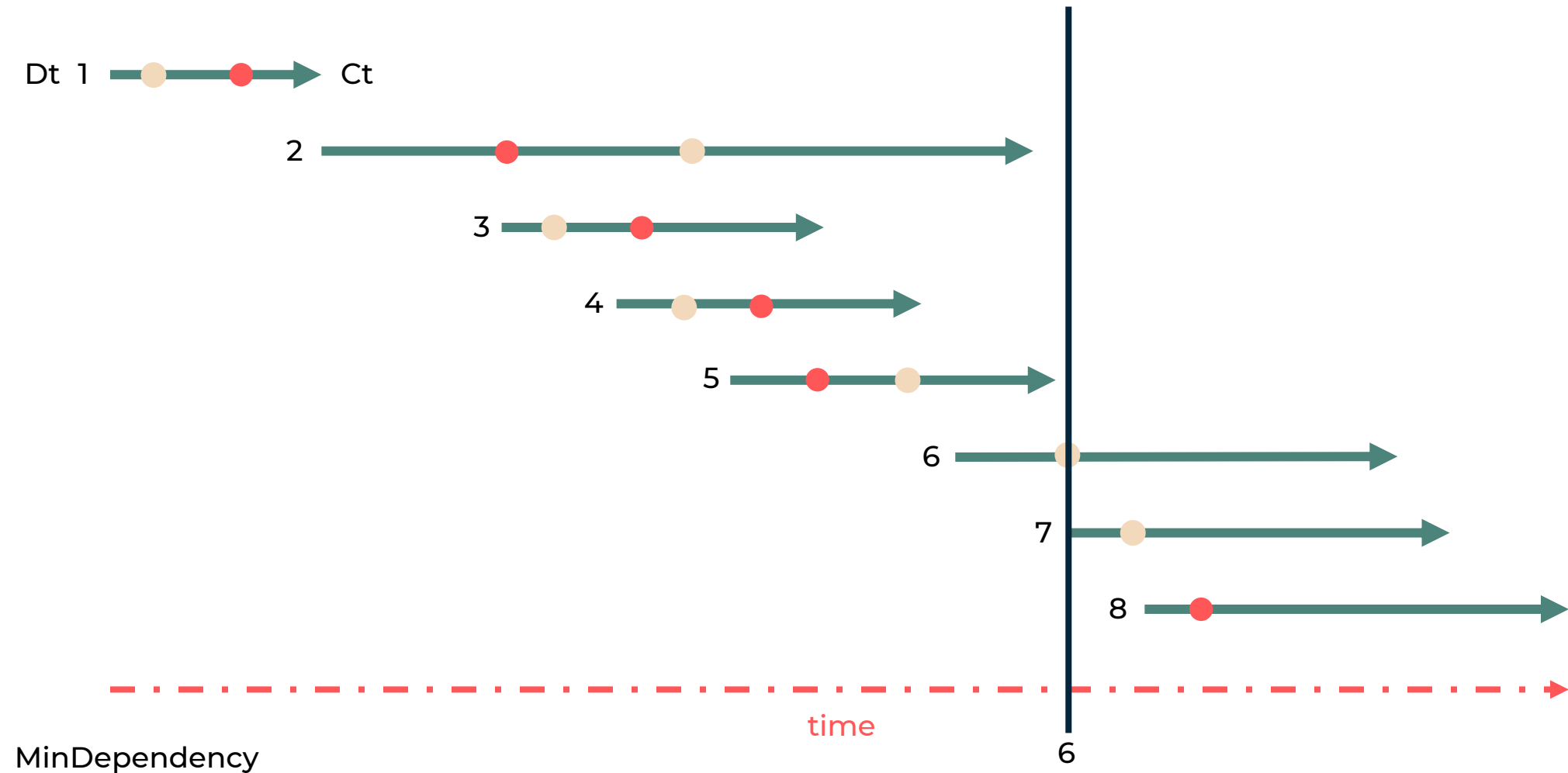
Bounded memory – Garbage collection



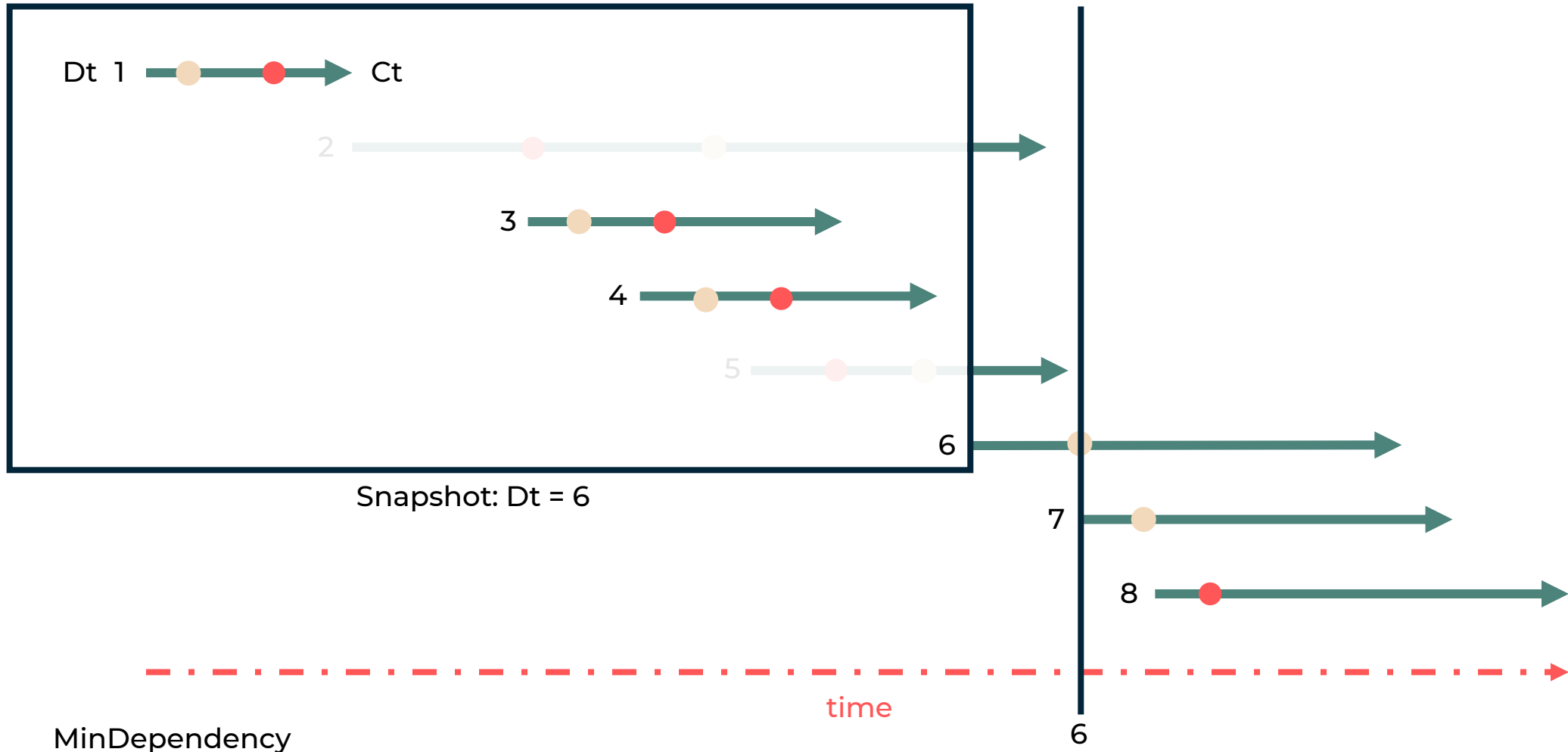
Bounded memory – Garbage collection



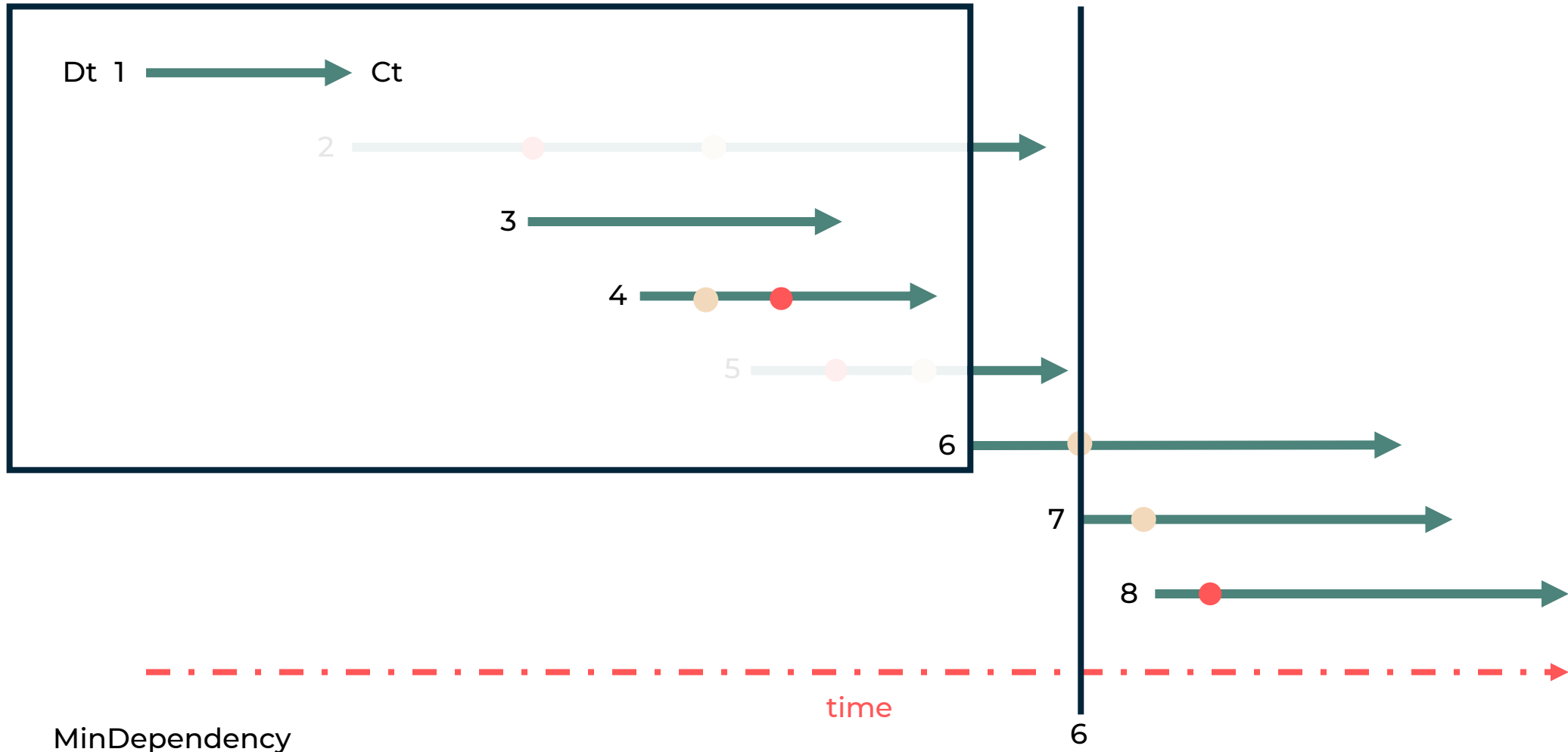
Bounded memory – Garbage collection



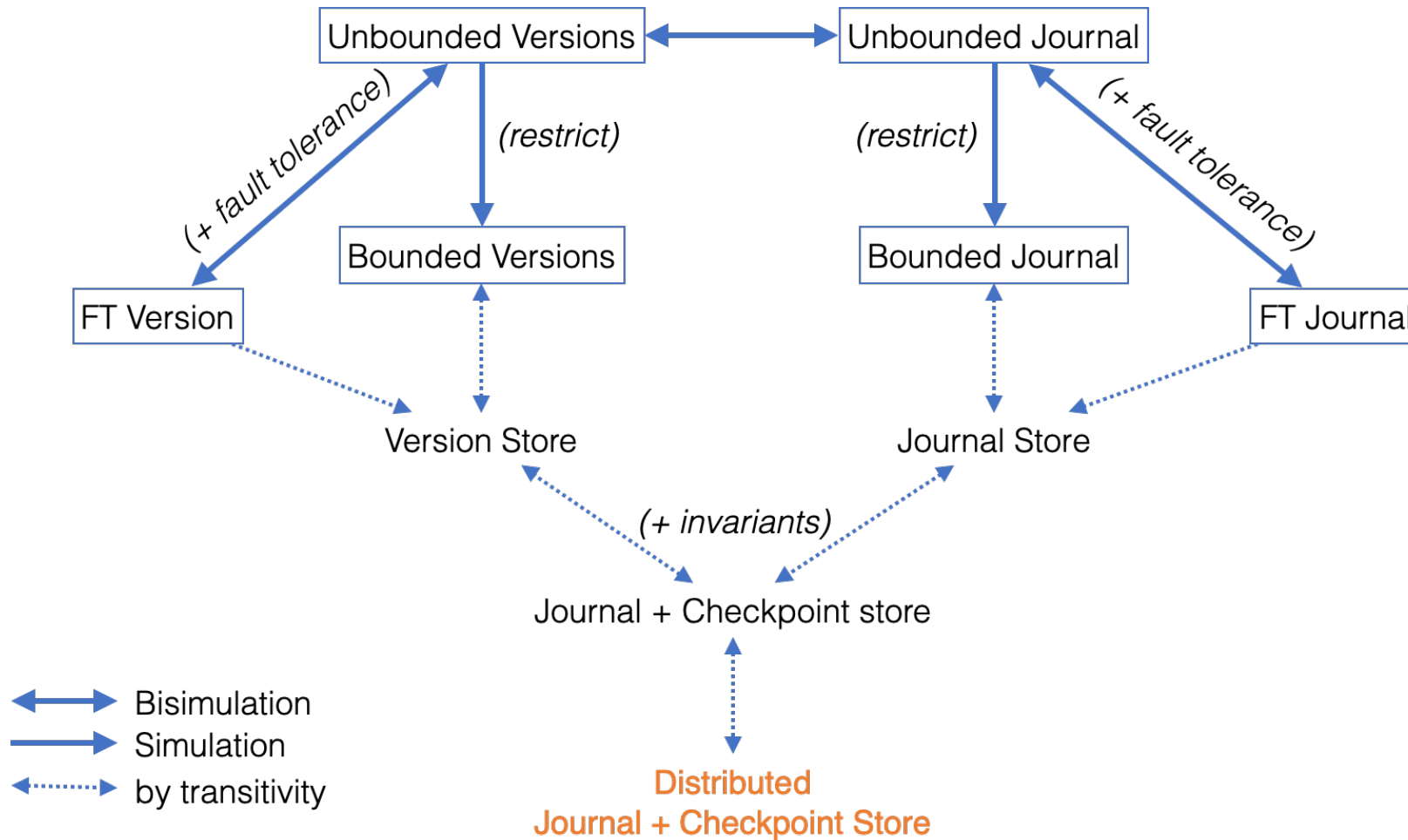
Bounded memory – Garbage collection



Bounded memory – Garbage collection



Roadmap



Ongoing and future work

Finishing all implementations

Testing invariants

Performance comparison

Static analysis