# SQLServerFast.com
## Execution Plan Video Training

Block 3: Combining data

Level: Advanced

Chapter 3: Hash Match (advanced)

# Hash Match

Hash Match
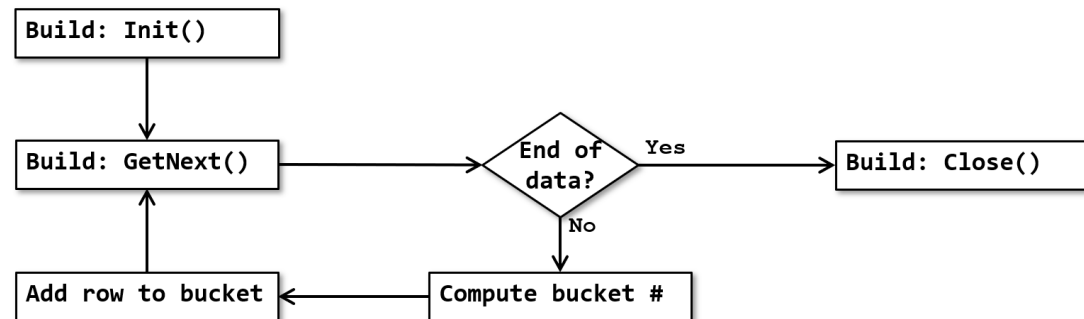Build phase
Probe phase


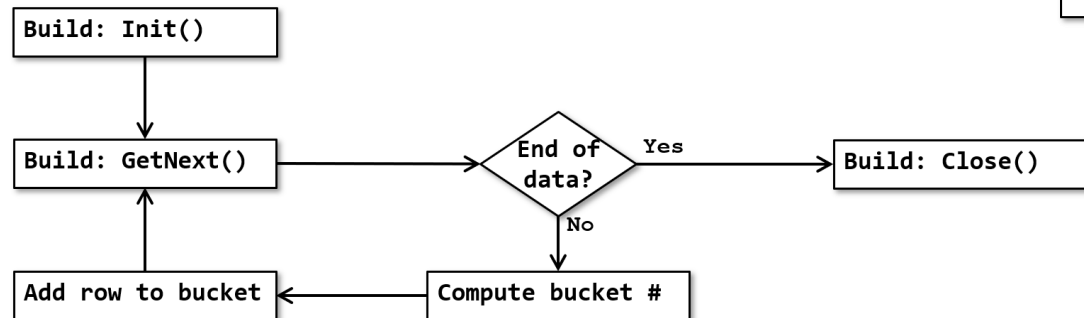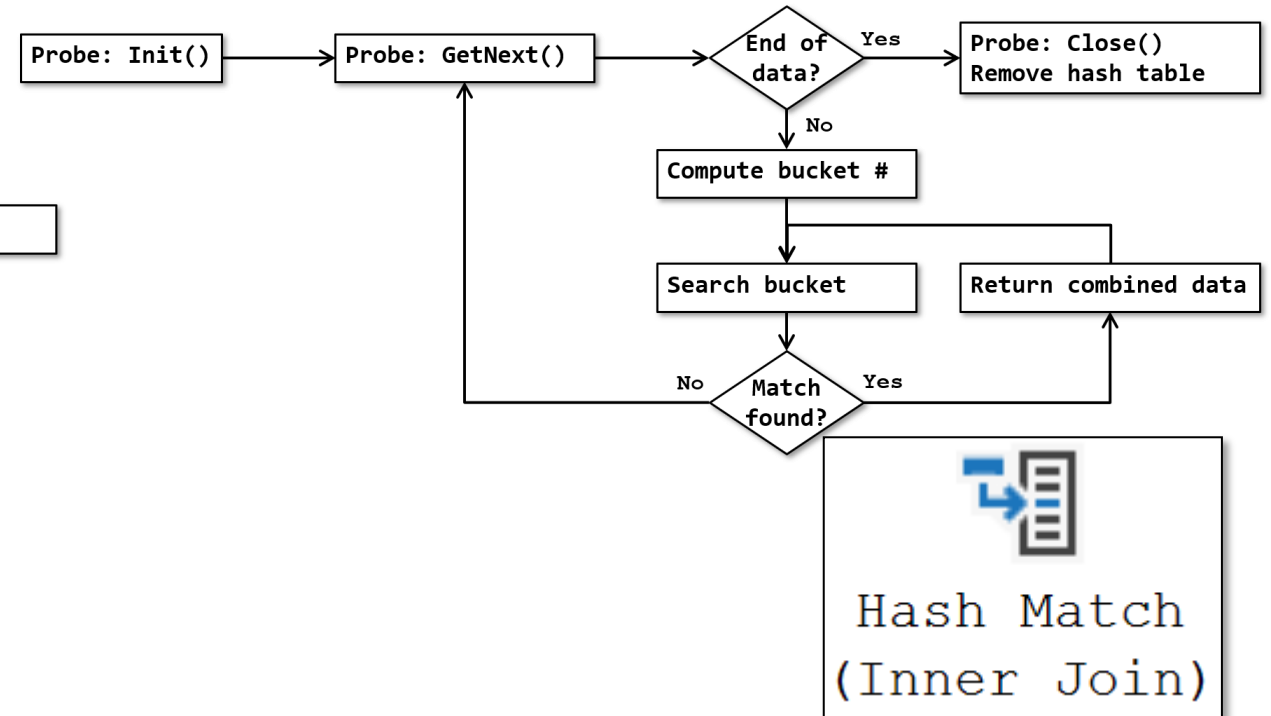
Hash Match
(Inner Join)

# Hash Match

## Hash Match
### Build phase

# Hash Match

## Hash Match

### Build phase

```
┌─────────────────┐
│ Build: Init()   │
└─────────────────┘
         │
         ▼
┌─────────────────┐              ╱╲
│ Build: GetNext()│───────────▶ ╱    ╲  Yes    ┌─────────────────┐
└─────────────────┘            ╱ End of╲──────▶│ Build: Close()  │
         ▲                     ╲ data? ╱       └─────────────────┘
         │                      ╲    ╱
         │                       ╲╱
         │                        │ No
         │                        ▼
┌─────────────────┐       ┌─────────────────┐
│Add row to bucket│◀──────│ Compute bucket #│
└─────────────────┘       └─────────────────┘
```

### Probe phase

```
┌─────────────┐      ┌──────────────────┐          ╱╲        Yes   ┌──────────────────┐
│ Probe: Init()│────▶│ Probe: GetNext() │─────────╱    ╲─────────▶ │ Probe: Close()   │
└─────────────┘      └──────────────────┘        ╱ End of╲         │ Remove hash table│
                              ▲                   ╲ data? ╱         └──────────────────┘
                              │                    ╲    ╱
                              │                     ╲╱
                              │                      │ No
                              │                      ▼
                              │             ┌──────────────────┐
                              │             │ Compute bucket # │
                              │             └──────────────────┘
                              │                      │
                              │          ┌───────────┴──────────┐
                              │          ▼                      │
                              │   ┌──────────────┐      ┌──────────────────────┐
                              │   │ Search bucket│      │ Return combined data │
                              │   └──────────────┘      └──────────────────────┘
                              │          │                      ▲
                              │         ╱╲                      │
                         No   │        ╱    ╲   Yes             │
                              └───────╱ Match ╲─────────────────┘
                                      ╲ found?╱
                                       ╲    ╱
                                        ╲╱
```

Hash Match
(Inner Join)

# Hash Match (inner to left outer)

From inner join to left outer join

Add unmatched left rows

These are from the build input

Can be found in the hash table

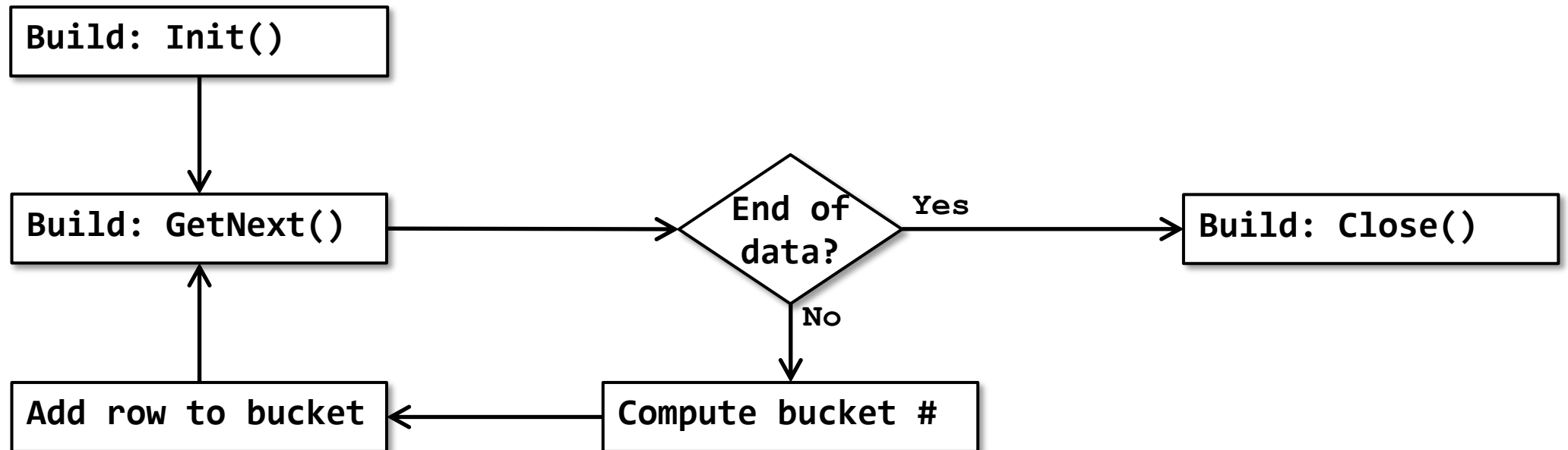But which of these rows to return?

Use a Boolean to track when there are matches

Hash Match
(Inner Join)

# Hash Match (inner to left outer)

From inner join to left outer join

Add unmatched left rows

These are from the build input

Can be found in the hash table

But which of these rows to return?

Use a Boolean *for each row in the hash table* to track when there are matches
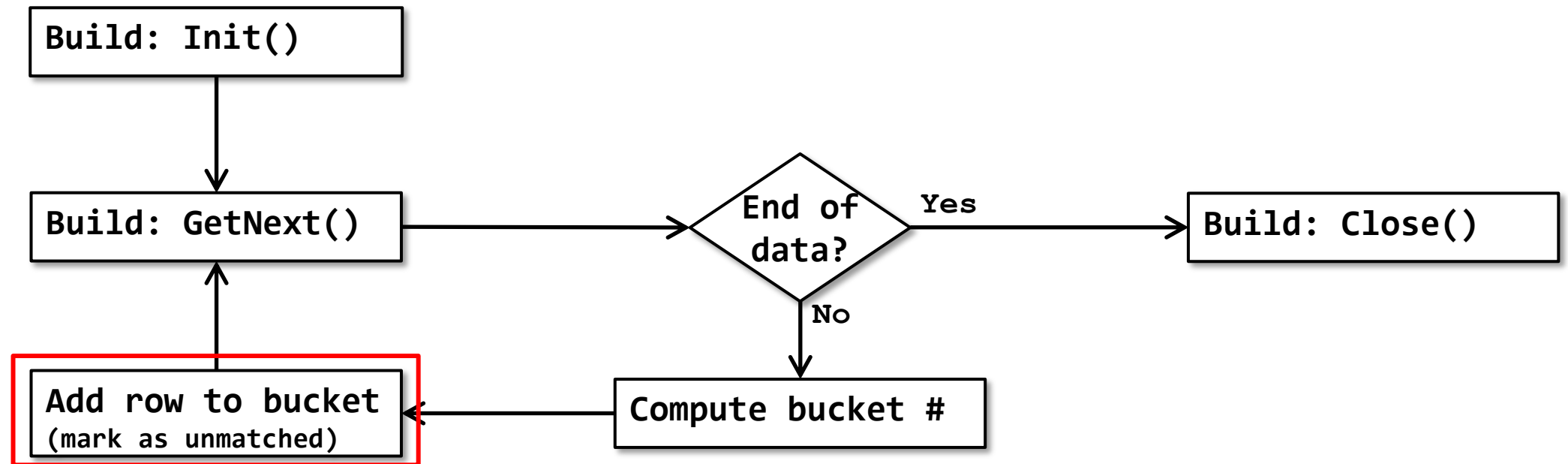


Hash Match
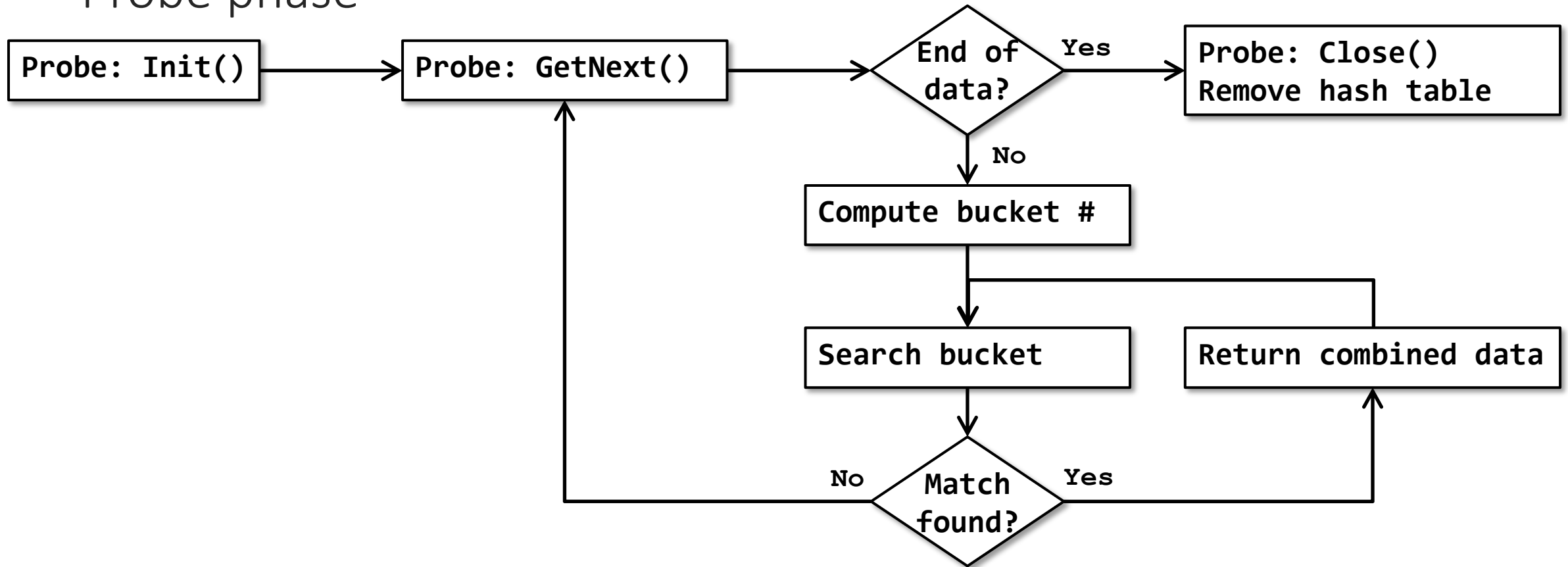(Inner Join)

# Hash Match (inner to left outer)

Build phase



```
Build: Init()
    |
    v
Build: GetNext()  ----->  End of data?  --Yes-->  Build: Close()
    ^                          |
    |                          No
    |                          |
    |                          v
Add row to bucket  <-----  Compute bucket #
```

# Hash Match (left outer join, no spill)

Build phase

# Hash Match (inner to left outer)

Probe phase

# Hash Match (inner to left outer)
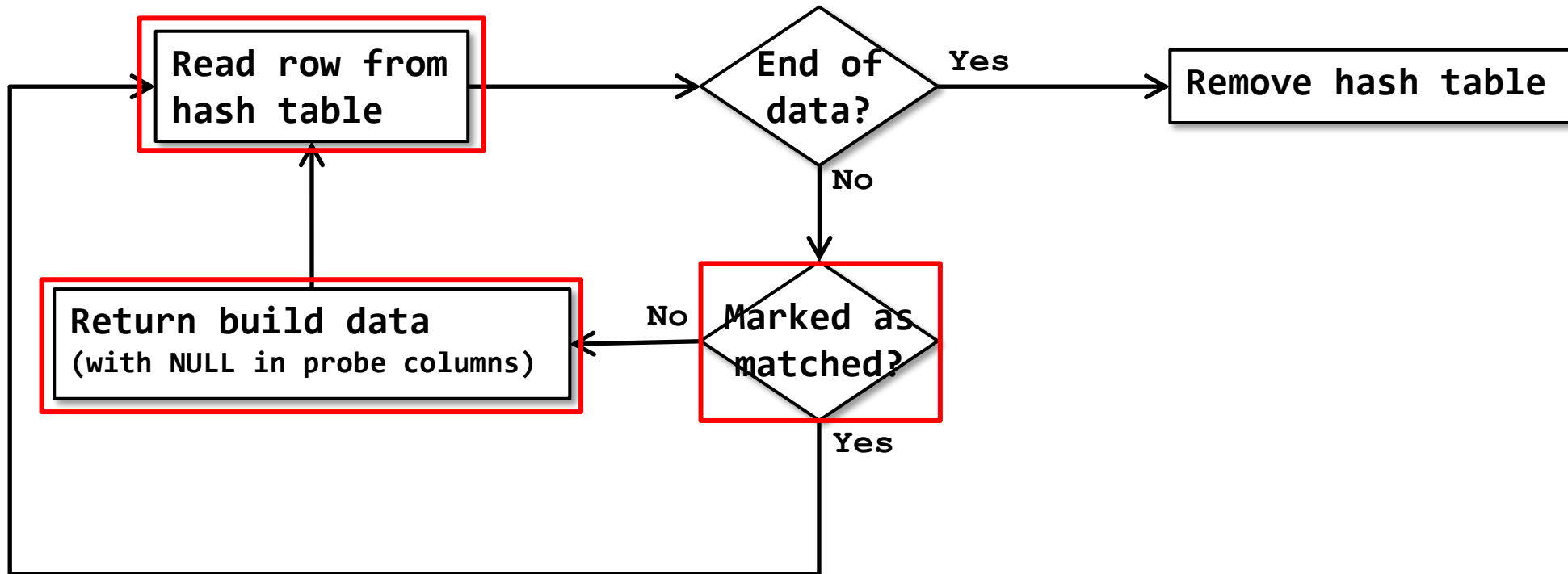
Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?  --Yes-->  Probe: Close()
                                                                Remove hash table

                                         |No
                                         ↓
                                   Compute bucket #
                                         ↓
                                   Search bucket        Return combined data
                                         ↓                       ↑
                                   Match found?  --Yes-->  Update row in bucket
                                         |No                (mark as matched)
```

# Hash Match (left outer join, no spills)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?
```

**End of data?** —Yes→ **Probe: Close() / Run final phase / Remove hash table**

**End of data?** —No→ **Compute bucket #**

**Compute bucket #** → **Search bucket** → **Match found?**

**Match found?** —No→ (back to Probe: GetNext())

**Match found?** —Yes→ **Update row in bucket (mark as matched)** → **Return combined data**

# Hash Match (left outer join, no spills)

Final phase (officially part of the probe phase)



```
Read row from          End of        Yes    Remove hash table
hash table             data?

                         No

                                     No    Marked as
Return build data                          matched?
(with NULL in probe columns)

                                            Yes
```

# Hash Match (inner join, no spills)

Probe phase

```
┌─────────────────┐        ┌─────────────────────┐         ◇ End of ◇  Yes   ┌──────────────────────┐
│ Probe: Init()   │───────▶│ Probe: GetNext()    │───────▶ ◇  data? ◇ ─────▶ │ Probe: Close()       │
└─────────────────┘        └─────────────────────┘         ◇         ◇        │ Remove hash table    │
                                                                 │ No          └──────────────────────┘
                                                                 ▼
                                                        ┌──────────────────┐
                                                        │ Compute bucket # │
                                                        └──────────────────┘
                                                                 │
                                                                 ▼
                                                        ┌──────────────────┐         ┌──────────────────────┐
                                                        │ Search bucket    │         │ Return combined data │
                                                        └──────────────────┘         └──────────────────────┘
                                                                 │
                                              No      ◇  Match  ◇   Yes
                                           ◀────────  ◇ found?  ◇  ─────▶
                                                      ◇         ◇
```

# Hash Match (inner to right outer)

Probe phase



```
Probe: Init()  →  Probe: GetNext()  →  End of data?  ──Yes──→  Probe: Close()
                                                                Remove hash table
                                              │
                                              No
                                              ↓
                                        Compute bucket #
                                              │
                                              ↓
                                        Search bucket          Return combined data
                                              │
                                              ↓
                            No ──  Match found?  ──Yes
```

*SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis*

# Hash Match (inner to right outer)

Probe phase

```
┌─────────────────┐       ┌─────────────────────────┐              ◇ End of ◇  Yes   ┌─────────────────────┐
│ Probe: Init()   │ ────▶ │ Probe: GetNext()        │ ────▶  ◇    data?    ◇ ────▶ │ Probe: Close()      │
└─────────────────┘       │ (mark as unmatched)     │              ◇         ◇        │ Remove hash table   │
                          └─────────────────────────┘                  │             └─────────────────────┘
                                                                       No
```

**Probe: Init()** → **Probe: GetNext()** (mark as unmatched) → **End of data?**

— Yes → **Probe: Close()** Remove hash table

— No ↓

**Compute bucket #**

↓

**Search bucket**

↓

**Match found?** — Yes → **Return combined data** (mark probe as matched)

— No → (back to Probe: GetNext())

# Hash Match (inner to right outer)

Probe phase

```
┌──────────────────┐        ┌─────────────────────────┐        ╱◆╲              ┌──────────────────────┐
│ Probe: Init()    │───────▶│ Probe: GetNext()        │──────▶ End of   ─Yes─▶ │ Probe: Close()       │
│                  │        │ (mark as unmatched)     │       data?            │ Remove hash table    │
└──────────────────┘        └─────────────────────────┘        ╲◆╱              └──────────────────────┘
```

End of data? → Yes → Probe: Close() Remove hash table

No

Compute bucket #

Search bucket

Return combined data
(mark probe as matched)

Probe row matched?

Yes — No

Match found?

No — Yes

# Hash Match (right outer join, no spills)

Probe phase



```
Probe: Init()  →  Probe: GetNext()
                   (mark as unmatched)
```

End of data? — Yes → Probe: Close()
Remove hash table

No ↓

Compute bucket #

Search bucket

Return combined data
(mark probe as matched)

Match found? — Yes →

Match found? — No → Probe row matched?

Probe row matched? — Yes →

Probe row matched? — No → Return probe data
(with NULL in build columns)

# Hash Match (right outer to full outer)

Probe phase

```
Probe: Init()
```

```
Probe: GetNext()
(mark as unmatched)
```

End of data? —— **Yes** →
```
Probe: Close()
Remove hash table
```

**No** ↓

```
Compute bucket #
```

```
Return probe data
(with NULL in build columns)
```

```
Search bucket
```

```
Return combined data
(mark probe as matched)
```

Probe row matched? —— **Yes**
**No** ↑

Match found? —— **No** ← —— **Yes** →

# Hash Match (right outer to full outer)

Probe phase



```
Probe: Init()
```

```
Probe: GetNext()
(mark as unmatched)
```

```
End of
data?
```
Yes

```
Probe: Close()
Remove hash table
```

No

```
Compute bucket #
```

```
Return probe data
(with NULL in build columns)
```

```
Search bucket
```

```
Return combined data
(mark probe as matched)
```

```
Probe row
matched?
```
No / Yes

```
Match
found?
```
No / Yes

```
Update row in bucket
(mark as matched)
```

*SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis*

# Hash Match (full outer join, no spills)

Probe phase

```
Probe: Init()
```

```
Probe: GetNext()
(mark as unmatched)
```

End of data?  — Yes →

```
Probe: Close()
Run final phase
Remove hash table
```

No ↓

```
Compute bucket #
```

```
Return probe data
(with NULL in build columns)
```

```
Search bucket
```

```
Return combined data
(mark probe as matched)
```

Probe row matched?
Yes / No

Match found?
No / Yes →

```
Update row in bucket
(mark as matched)
```

SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis

# Hash Match (full outer join, no spills)

Typical output order *(when there are no spills!)*

    All probe rows, matched and unmatched, in original order

    All unmatched build rows, in "semi random" order

This order is not guaranteed

    But you should still be aware of it, in case it matters



Hash Match
(Inner Join)

# Hash Match (inner to left semi)

From inner join to left semi join

    Only return matched build rows

    Return matched build rows once, regardless of number of matches

    Two options

        Return when first match found

            Mark as matched

            Skip on next match

        Mark as matched when any match found

            Return matched rows in final phase



Hash Match
(Inner Join)

# Hash Match (left semi join, no spills)

From inner join to left semi join

 Only return matched build rows

 Return matched build rows once, regardless of number of matches

 Marks as matched during probe phase

 Returns matched rows during final phase

 Operator is now *fully* blocking

  No data returned during build phase

  No data returned during probe phase

  All data returned during final phase

Hash Match
(Inner Join)

# Hash Match (left outer to left semi)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?
```

End of data? — Yes → **Probe: Close() / Run final phase / Remove hash table**

End of data? — No ↓

**Compute bucket #**

↓

**Search bucket** → **Return combined data**

↓

**Match found?** — No → (back to Probe: GetNext())

**Match found?** — Yes → **Update row in bucket (mark as matched)** → **Return combined data**

# Hash Match (left semi join, no spills)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?
```

End of data? —Yes→ Probe: Close() / Run final phase / Remove hash table

End of data? —No→ Compute bucket # → Search bucket → Match found?

Match found? —No→ (back to Probe: GetNext())

Match found? —Yes→ Update row in bucket (mark as matched)

# Hash Match (left semi join, no spills)

Final phase

# Hash Match (left anti semi join, no spills)

Probe phase



```
Probe: Init()  →  Probe: GetNext()  →  End of data?
                                         Yes →  Probe: Close()
                                                Run final phase
                                                Remove hash table
                                         No ↓
                                       Compute bucket #
                                              ↓
                                       Search bucket
                                              ↓
                                       Match found?
                            No ←                    → Yes → Update row in bucket
                                                                  (mark as matched)
```

# Hash Match (left anti semi join, no spills)

Final phase



SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis

# Hash Match (probed left semi join)

Probed left semi join

    Not supported

        Not clear why

        Would be relatively easy to build

        Never encountered, never been able to repro



Hash Match
(Inner Join)

# Hash Match (inner to right semi)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?  --Yes-->  Probe: Close()
                                                                 Remove hash table
```

End of data? --No--> Compute bucket #

Compute bucket # → Search bucket → Match found?

Match found? --Yes--> Return combined data

Match found? --No--> (back to Probe: GetNext())

Return combined data (from Compute bucket # path)

# Hash Match (inner to right semi)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  ◇ End of data? ─Yes→  Probe: Close() / Remove hash table
                                              │
                                              No
                                              ↓
                                        Compute bucket #
                                              │
                                              ↓
                                        Search bucket        Return probe data
                                              │
                                              ↓
                                        ◇ Match found? ─Yes→
                                              │
                                              No
```

# Hash Match (right semi join, no spills)

Probe phase



```
┌──────────────┐        ┌──────────────────┐        ◇ End of ◇  Yes   ┌──────────────────┐
│ Probe: Init()│───────▶│ Probe: GetNext() │──────▶◇   data?  ◇──────▶│ Probe: Close()   │
└──────────────┘        └──────────────────┘        ◇         ◇       │ Remove hash table│
                                                         │ No         └──────────────────┘
                                                         ▼
                                              ┌──────────────────┐
                                              │ Compute bucket # │
                                              └──────────────────┘
                                                         │
                                                         ▼
                                              ┌──────────────────┐   ┌──────────────────┐
                                              │  Search bucket   │   │ Return probe data│
                                              └──────────────────┘   └──────────────────┘
                                                         │
                                                         ▼
                                    No       ◇  Match   ◇   Yes
                                  ◀──────────◇  found?  ◇──────────▶
                                             ◇          ◇
```

# Hash Match (right semi to right anti semi)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  ◇ End of data?
```

◇ End of data? —Yes→ **Probe: Close() Remove hash table**

End of data? —No↓

**Compute bucket #**

↓

**Search bucket**

↓

◇ **Match found?** —Yes→ **Return probe data**

Match found? —No→ (back to Probe: GetNext())

# Hash Match (right semi to right anti semi)

Probe phase

Probe: Init() → Probe: GetNext() → End of data? —Yes→ Probe: Close() Remove hash table

End of data? —No→ Compute bucket # → Search bucket → Match found?

Match found? —No→ (loops back to GetNext)

Match found? —Yes→ (loops back)

[red empty box]

# Hash Match (right anti semi join, no spills)

Probe phase

```
┌─────────────────┐        ┌──────────────────┐         ◇ End of        Yes    ┌──────────────────────┐
│ Probe: Init()   │───────▶│ Probe: GetNext() │────────▶◇ data?    ─────────▶ │ Probe: Close()       │
└─────────────────┘        └──────────────────┘         ◇               │     │ Remove hash table    │
                                                              │ No             └──────────────────────┘
                                                              ▼
                                                   ┌──────────────────┐
                                                   │ Compute bucket # │
                                                   └──────────────────┘
                           ┌──────────────────┐            │
                           │ Return probe data│            ▼
                           └──────────────────┘   ┌──────────────────┐
                                                   │ Search bucket    │
                                                   └──────────────────┘
                                                            │
                                          No        ◇ Match      Yes
                                     ─────────────  ◇ found?  ──────────
```

# Hash Match (right anti semi join, no spills)

Probe phase (alternative version)

```
Probe: Init()
```

```
Probe: GetNext()
```

End of data?

**Yes** → Probe: Close() Remove hash table

**No** ↓

```
Compute bucket #
```

```
Return probe data
```

```
Search bucket
```

(mark probe as matched)

Probe row matched?

**No**

**Yes**

Match found?

**No**

**Yes**

# Hash Match (union)

Union

   Merge Join (Union)

      Both inputs have to be free of duplicates already

      Avoids new duplicates after combining the two inputs

   Hash Match (Union)

      Probe input has to be free of duplicates already

      Duplicates in build input are removed by the operator

      Avoids new duplicates after combining the two inputs



Hash Match
(Inner Join)

# Hash Match (inner join to union)

Probe phase

```
┌─────────────────┐        ┌─────────────────┐                    ┌──────────────────────┐
│ Probe: Init()   │───────▶│ Probe: GetNext()│──────▶  End of ──Yes──▶│ Probe: Close()       │
└─────────────────┘        └─────────────────┘         data?           │ Remove hash table    │
                                                                       └──────────────────────┘
```

End of data?

Yes

Probe: Close()
Remove hash table

No

Compute bucket #

Search bucket

Return combined data

Match found?

No

Yes

# Hash Match (inner join to union)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?  --Yes-->  Probe: Close()
                                                                 Remove hash table
```

End of data? --No--> Compute bucket # --> Search bucket --> Match found?

Match found? --No--> (back to Probe: GetNext())

Match found? --Yes--> (red box)

# Hash Match (inner join to union)

Probe phase

# Hash Match (union, no spills)

Probe phase

```
┌──────────────────┐        ┌──────────────────┐         ◇ End of ◇   Yes    ┌──────────────────────┐
│ Probe: Init()    │───────▶│ Probe: GetNext() │───────▶◇  data?  ◇────────▶│ Probe: Close()       │
└──────────────────┘        └──────────────────┘         ◇        ◇         │ Remove hash table    │
                                                              │              └──────────────────────┘
                                                            No│
                                                              ▼
                                               ┌──────────────────────┐
                                               │ Compute bucket #     │
                                               └──────────────────────┘
                            ┌──────────────────┐              │
                            │ Return probe data│              ▼
                            └──────────────────┘   ┌──────────────────────┐
                                                   │ Search bucket        │
                                                   └──────────────────────┘
                                                              │
                                                              ▼
                                            No            ◇ Match ◇    Yes
                                        ◀──────────────◇  found? ◇──────────────▶
                                                        ◇        ◇
```

# Hash Match (union, no spills)

Probe phase (alternative version)

```
Probe: Init()  →  Probe: GetNext()  →  End of data? --Yes-->  Probe: Close()
                                                               Remove hash table
```

End of data? --No--> Compute bucket #

```
Return probe data
```

```
Compute bucket #  →  Search bucket  →  Match found? --Yes-->  (mark probe as matched)
```

Match found? --No--> Probe row matched?

Probe row matched? --Yes--> (loop back to GetNext)

Probe row matched? --No--> Return probe data

# Hash Match (union, no spills)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?
```

End of data? — Yes → **Probe: Close() / Run final phase / Remove hash table**

End of data? — No ↓

**Compute bucket #**

↓

**Search bucket**

↓

**Match found?** — No → (return to GetNext loop)

**Match found?** — Yes → (loop back)

**Return probe data** → Probe: GetNext()

# Hash Match (union, no spills)

Final phase

# Hash Match (union, no spills)

Final phase



```
┌─────────────────┐                    ◇ End of ◇          ┌──────────────────┐
│ Read row from   │ ──────────────────▶  data?   ── Yes ──▶│ Remove hash table │
│ hash table      │                    ◇         ◇         └──────────────────┘
└─────────────────┘                        │
        ▲                                  No
        │                                  │
┌─────────────────┐                        │
│ Return build data│◀──────────────────────┘
└─────────────────┘
```

# Hash Match (union, no spills)

Final phase (bad alternative, not used!!!)

# Hash Match (inner to union)

Build phase



```
┌─────────────────────┐
│ Build: Init()       │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐          ◇ End of      Yes    ┌──────────────────┐
│ Build: GetNext()    │ ────────▶  data?      ──────▶ │ Build: Close()   │
└─────────────────────┘          ◇                    └──────────────────┘
          ▲                         │ No
          │                         ▼
┌─────────────────────┐          ┌──────────────────┐
│ Add row to bucket   │ ◀─────── │ Compute bucket # │
└─────────────────────┘          └──────────────────┘
```

# Hash Match (inner to union)

Build phase

# Hash Match (union, no spills)

Final phase

# Hash Match (union, no spills)

Build phase

# Hash Match (union, no spills)

Probe phase

```
Probe: Init()  →  Probe: GetNext()  →  End of data?  --Yes→  Probe: Close()
                                                                Run final phase
                                                                Remove hash table

                        Return probe data              │No

                                                   Compute bucket #

                                                   Search bucket

                              No ←  Match found?  Yes→
```

# Hash Match (all operations, no spills)

Build phase

```
┌─────────────────────┐
│ Build: Init()       │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐        ╱◇╲              ┌─────────────────────┐
│ Build: GetNext()    │─────▶ ◇ End of ◇ ─Yes─▶ │ Build: Close()      │
└─────────────────────┘        ◇ data? ◇        └─────────────────────┘
          ▲                     ╲◇╱
          │                      │ No
          │                      ▼
┌─────────────────────┐  ┌──────┐  ┌─────────────────────┐
│ Add row to bucket   │◀─│ All  │  │ Compute bucket #    │
│ (mark as unmatched) │  │joins │◀─┤                     │
└─────────────────────┘  └──────┘  └─────────────────────┘
          ▲                              │ Union
          │                              ▼
          │                            ╱◇╲
          │                     No ── ◇ Match ◇
          └───────────────────────── ◇ found? ◇
                                       ╲◇╱
                                        │ Yes
```

*SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis*

# Hash Match (all operations, no spills)

Probe phase



SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis

# Hash Match (all operations, no spills)

Final phase

# Hash Match

Memory

Hash table stored in memory

Memory Grant

Determined by optimizer

Based on estimates

May be adjusted for later executions by Memory Grant Feedback

Since SQL Server 2017 for batch mode

Since SQL Server 2019 for row mode

No additional memory allocations once query runs

(there are some rare exceptions to this in batch mode plans only)



Hash Match
(Inner Join)

# Hash Match

Memory

    Hash table stored in memory

    Memory Grant

    What if the build input is larger than the available memory?

        Fail with run-time error

        Different version of algorithm

            Still returns correct results

            Performance suffers



Hash Match
(Inner Join)

# Hash Match

Memory

  Hash table stored in memory

  Memory Grant

  What if the build input is larger than the available memory?

   ~~Fail with run-time error~~

   Different version of algorithm

     Still returns correct results

     Performance suffers

     Uses tempdb to store data that doesn't fit in memory

       This data is then read back and processed later

       Called "spilling" to tempdb

       Indicated in execution plan plus run-time statistics



Hash Match
(Inner Join)

# Hash Match

Memory

    Hash table stored in memory

    Memory Grant

    What if the build input is larger than the available memory?

       ~~Fail with run-time error~~

       Different version of algorithm

          Still returns correct results

          Performance suffers

          Uses tempdb to store data that doesn't fit in memory

             This data is then read back and processed later

             Called "spilling" to tempdb

             Indicated in execution plan plus run-time statistics, Extended Events, and SQL Trace



Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

Can only occur during the build phase

Build row read, but no memory available to store it

Algorithm changes from "in-memory hash join" to "dynamic hash join" or "grace hash join"



Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

Dynamic hash join / grace hash join

Hash table divided into several partitions

Partition number determined by hashing the bucket number

One "active" partition remains in memory

(For grace hash join, zero partitions remain in memory)

All inactive partitions spill to tempdb

Data already in hash table moved to "files" in tempdb

Rest of build phase stores data in memory or in appropriate file

End of build phase: active partition in memory, rest in files

Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

Dynamic hash join / grace hash join

Build phase: Active partition in memory, inactive partitions in tempdb files

Probe phase:

Rows in active partition can be regularly processed

Rows in inactive partitions stored in yet more tempdb files

This finds matches and non-matches for active partition only

Execute final phase (if needed) to complete final results for active partition

No results at all for inactive partitions

But build and probe data for these partitions is now in separate files in tempdb

Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

Dynamic hash join / grace hash join

Build phase: Active partition in memory, inactive partitions in tempdb files

Probe phase: Results for active partition, inactive partitions in tempdb files

Process inactive partitions

Make active

Read build file, store data in hash table

Read probe file, handle matches and non-matches as they are found

Execute final phase (if needed)

Repeat for each inactive partition



Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

Dynamic hash join / grace hash join

Multiple iterations of build, probe, and final phase

First iteration

Reads build input, stores in memory or in tempdb

Reads probe input, produces partial results or stores in tempdb

Later iterations

Reads build data from tempdb, stores in memory

Reads probe data from tempdb, produces partial results

Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

    Dynamic hash join / grace hash join

    Recursive hash join

        Initial partitions too large

        Partitions are split into new, smaller partitions



Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

    Dynamic hash join / grace hash join

    Recursive hash join

    Bail-out

        When recursive hash join fails

        Merge Join or Nested Loops for only those partition(s)

            Input from files in tempdb

Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

    Dynamic hash join / grace hash join

    Recursive hash join

    Bail-out

    Bit-vector filtering

        Probe rows that match empty bucket are handled immediately

        This reduces size of probe files for inactive partitions

Hash Match
(Inner Join)

# Hash Match

Hash spill in detail

    Dynamic hash join / grace hash join

    Recursive hash join

    Bail-out

    Bit-vector filtering

    Dynamic role reversal

        Partitions that start inactive

        Reverse build and probe if probe file has less rows



Hash Match
(Inner Join)

# Hash Match

Hash spill
  Data split into multiple partitions
  Data for partitions written to tempdb
  Partitions processed one by one

Output
  Start with first partition
    (may be mixed with some unmatched probe rows)
  Then second partition, etc
  Effectively (semi) random

Hash Match
(Inner Join)

# Hash Match

Multiple Hash Match operators in a single execution plan

Reusing memory

Hash teams



Hash Match
(Inner Join)

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory

# Hash Match

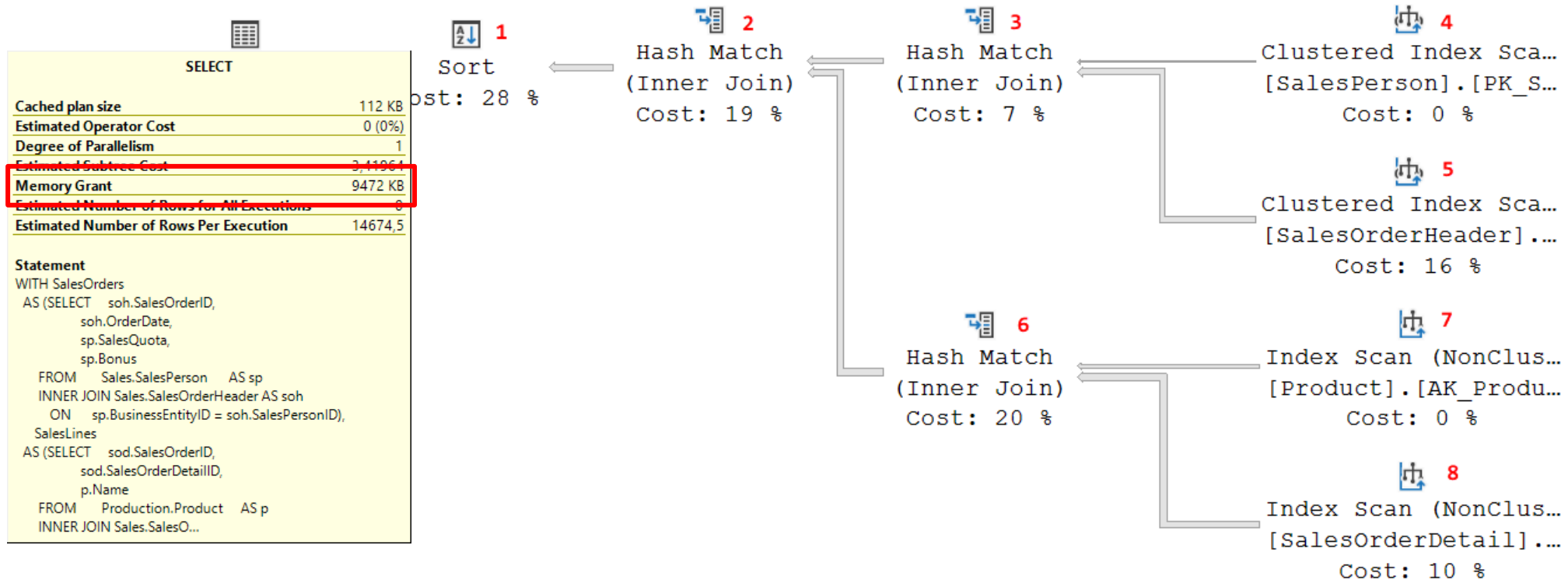Multiple *memory using* operators in a single execution plan

Reusing memory



*SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis*
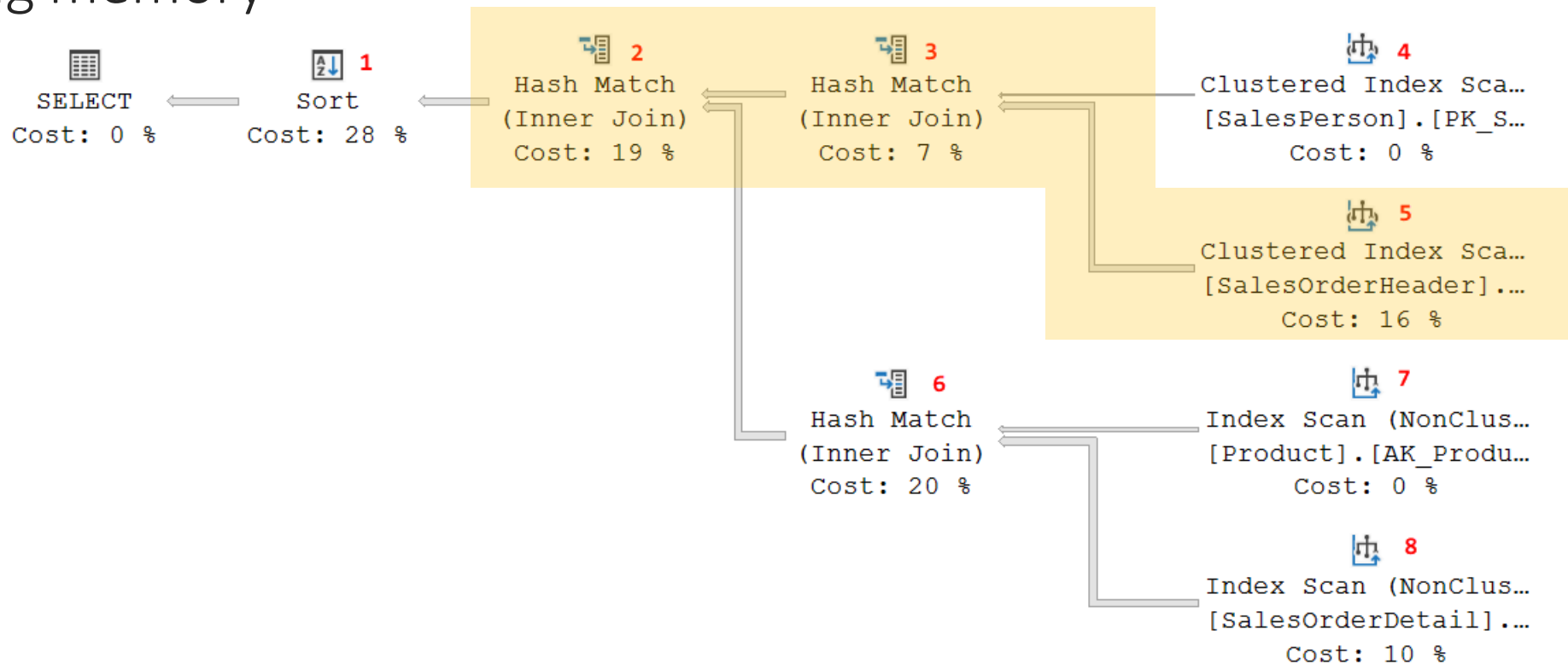
# Hash Match

Multiple *memory using* operators in a single execution plan
Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan
Reusing memory
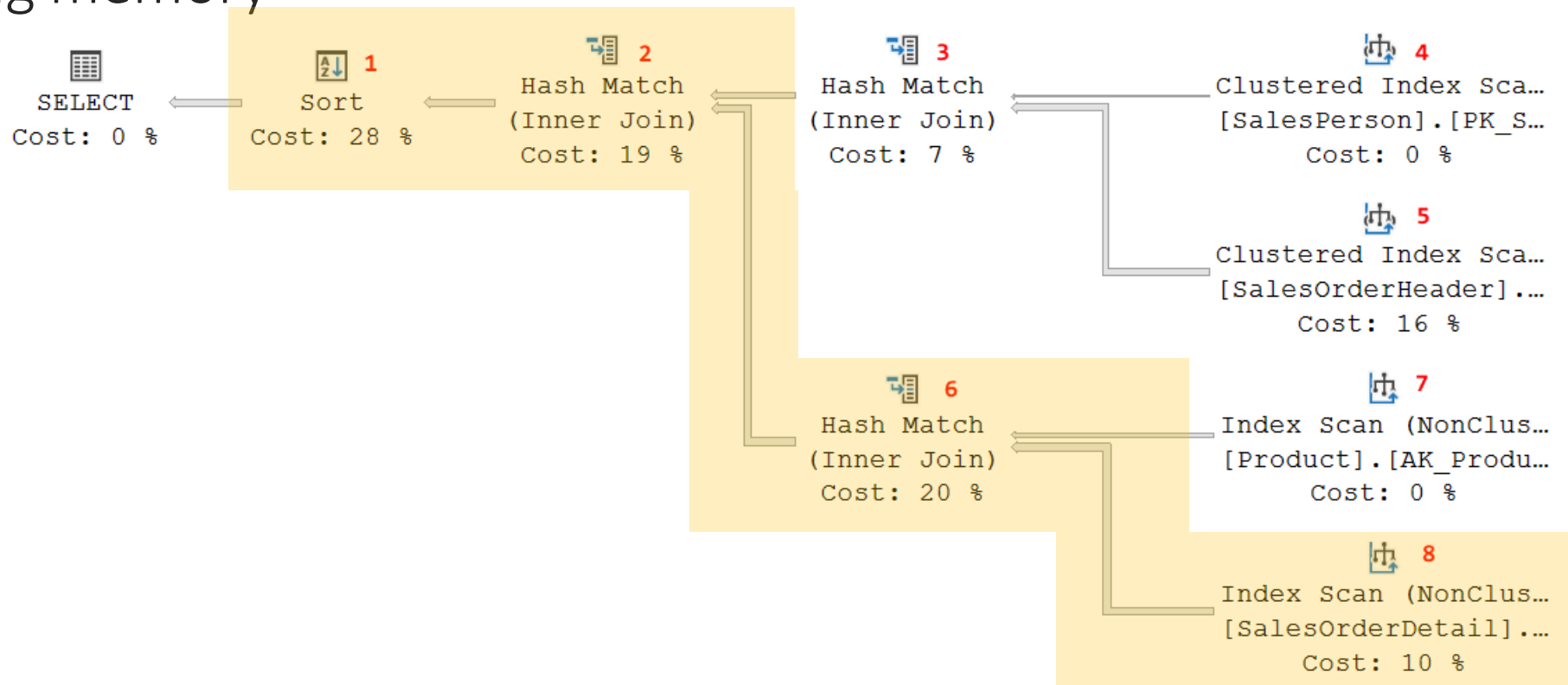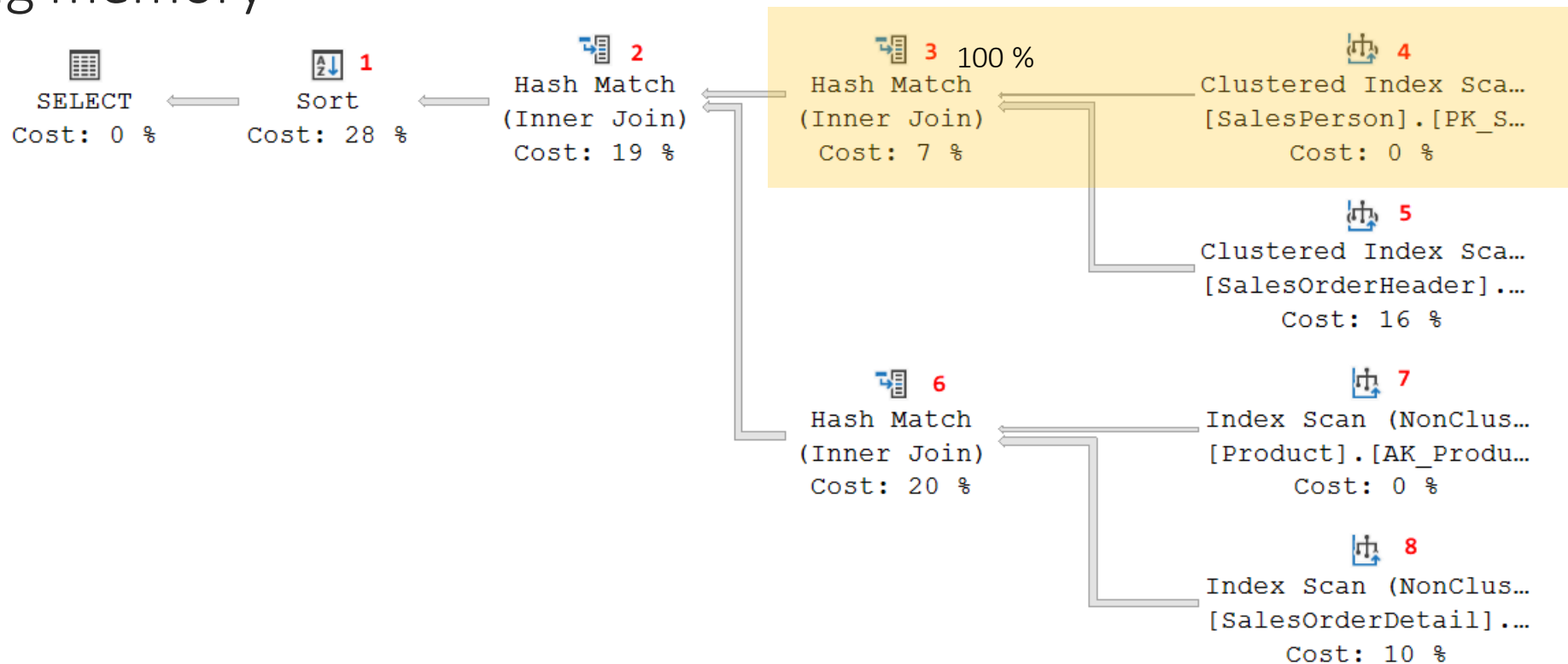
# Hash Match

Multiple *memory using* operators in a single execution plan
Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory



| | | |
|---|---|---|
| SELECT<br>Cost: 0 % | Sort   **1**<br>Cost: 28 % | |

Hash Match   **2**<br>(Inner Join)<br>Cost: 19 %

Hash Match   **3**<br>(Inner Join)<br>Cost: 7 %

Clustered Index Sca…   **4**<br>[SalesPerson].[PK_S…<br>Cost: 0 %

Clustered Index Sca…   **5**<br>[SalesOrderHeader].…<br>Cost: 16 %

Hash Match   **6**<br>(Inner Join)<br>Cost: 20 %

Index Scan (NonClus…   **7**<br>[Product].[AK_Produ…<br>Cost: 0 %

Index Scan (NonClus…   **8**<br>[SalesOrderDetail].…<br>Cost: 10 %

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory



| | |
|---|---|
| **SELECT** | |
| **Cached plan size** | 112 KB |
| **Estimated Operator Cost** | 0 (0%) |
| **Degree of Parallelism** | 1 |
| Estimated Subtree Cost | 3.4196... |
| **Memory Grant** | 9472 KB |
| Estimated Number of Rows for All Executions | 0 |
| **Estimated Number of Rows Per Execution** | 14674,5 |

**Statement**
WITH SalesOrders
 AS (SELECT   soh.SalesOrderID,
         soh.OrderDate,
         sp.SalesQuota,
         sp.Bonus
    FROM    Sales.SalesPerson    AS sp
    INNER JOIN Sales.SalesOrderHeader AS soh
      ON   sp.BusinessEntityID = soh.SalesPersonID),
    SalesLines
 AS (SELECT   sod.SalesOrderID,
         sod.SalesOrderDetailID,
         p.Name
    FROM    Production.Product   AS p
    INNER JOIN Sales.SalesO...

**1** Sort
Cost: 28 %

**2** Hash Match (Inner Join)
Cost: 19 %

**3** Hash Match (Inner Join)
Cost: 7 %

**4** Clustered Index Sca...
[SalesPerson].[PK_S...
Cost: 0 %

**5** Clustered Index Sca...
[SalesOrderHeader]....
Cost: 16 %

**6** Hash Match (Inner Join)
Cost: 20 %

**7** Index Scan (NonClus...
[Product].[AK_Produ...
Cost: 0 %

**8** Index Scan (NonClus...
[SalesOrderDetail]....
Cost: 10 %

*SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis*
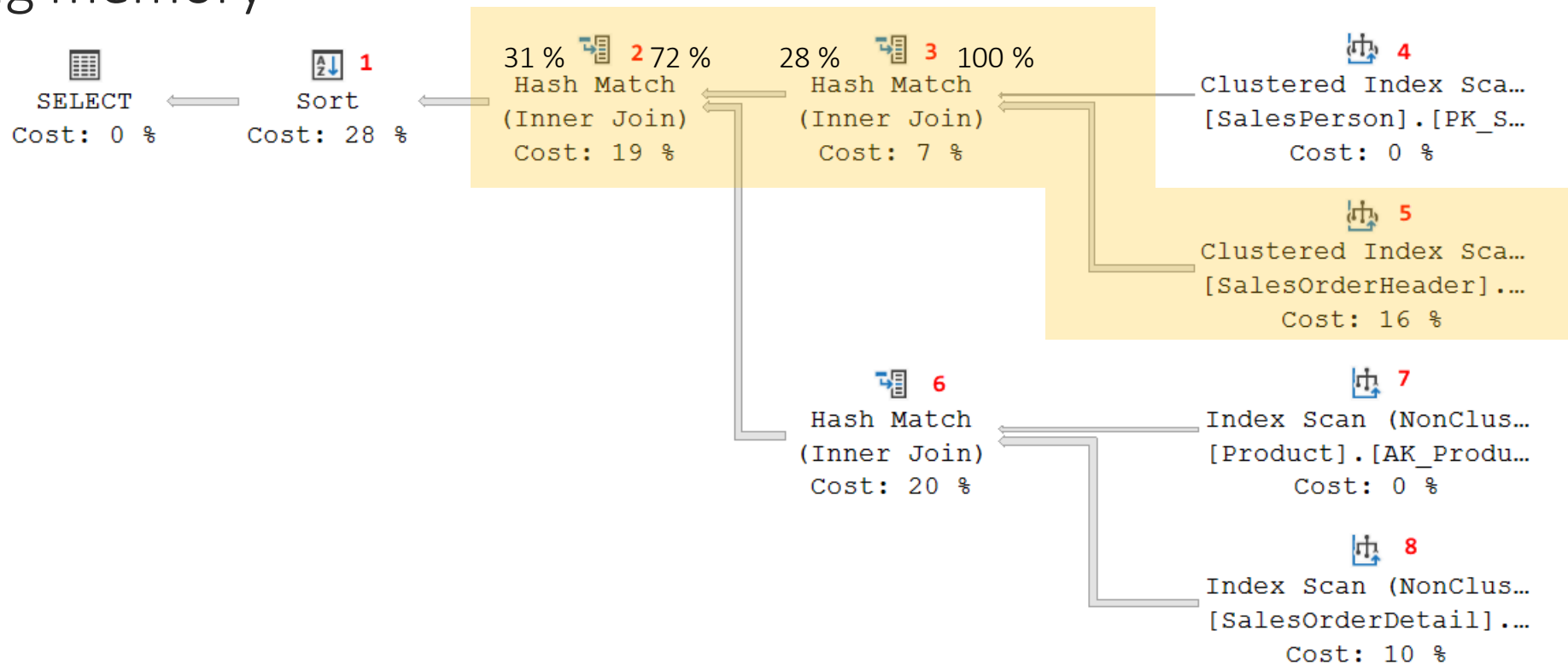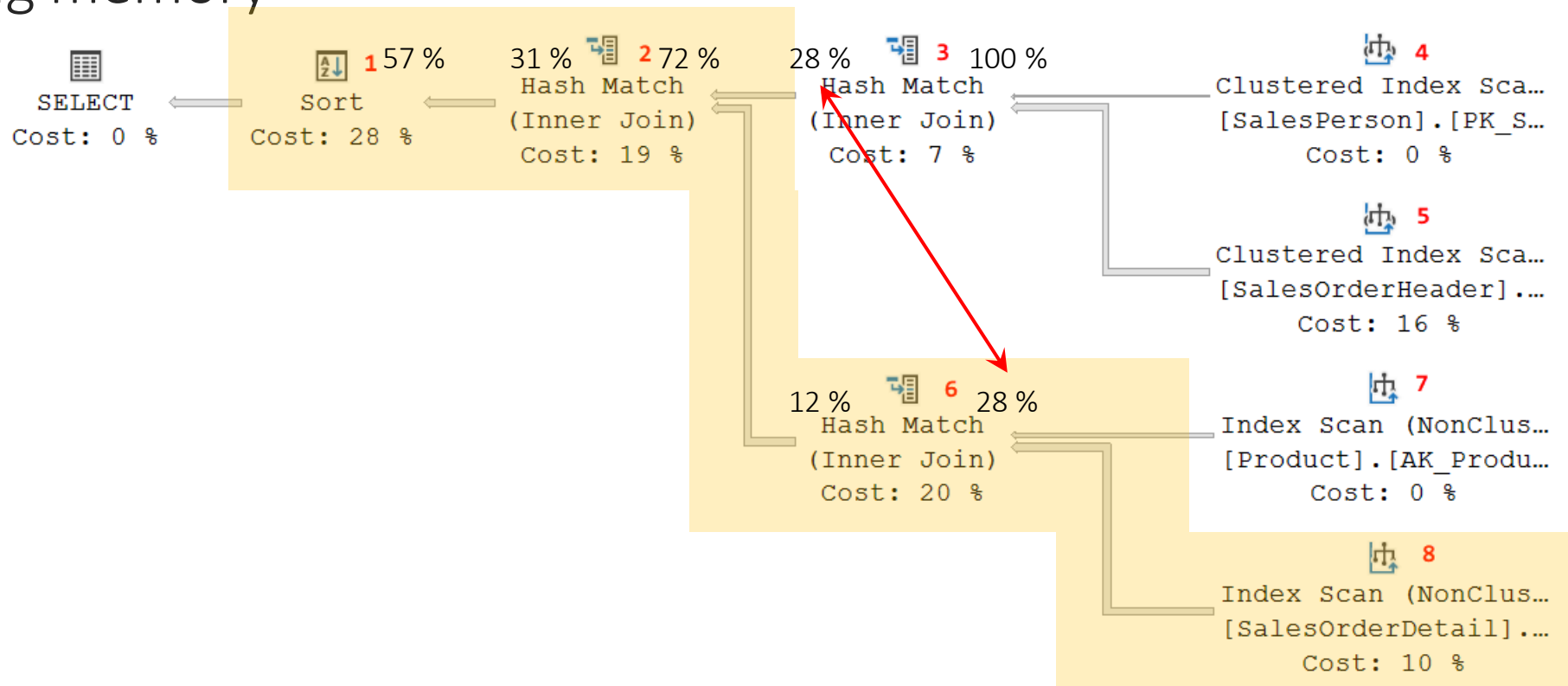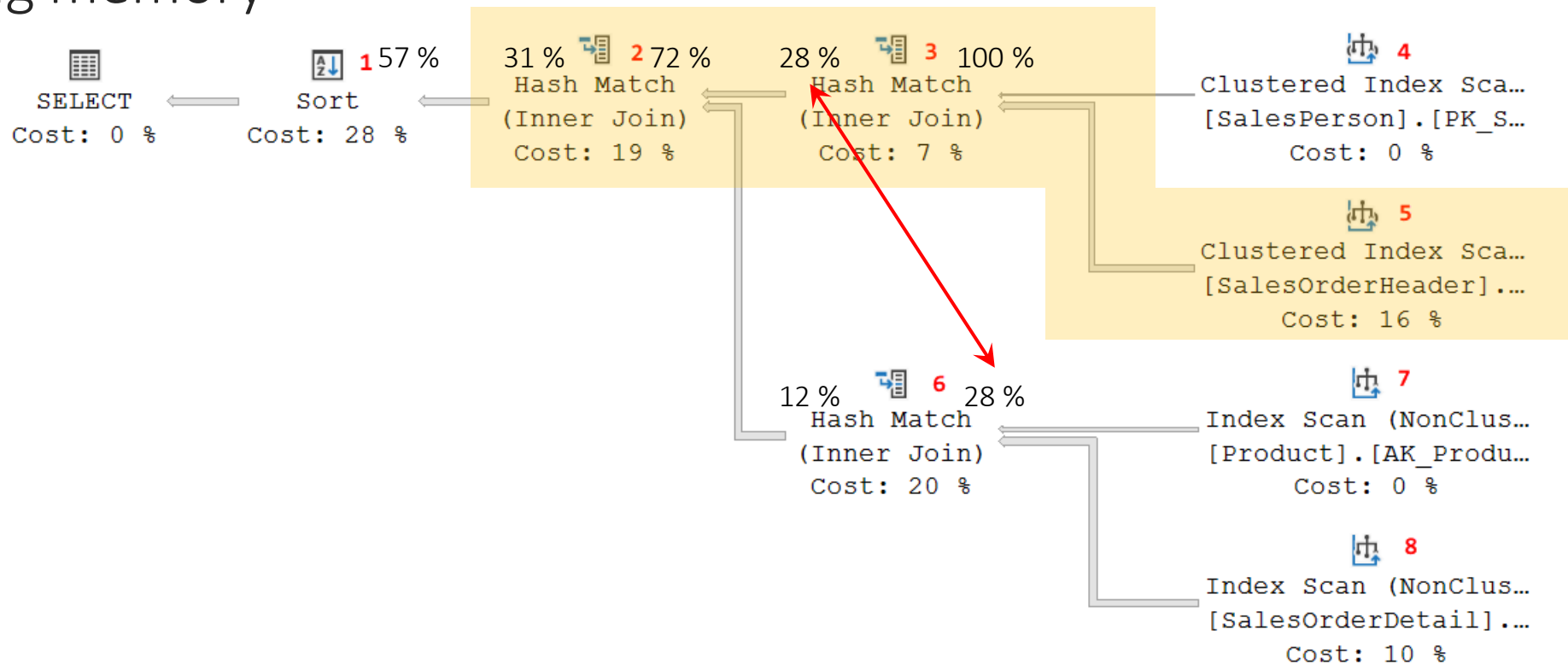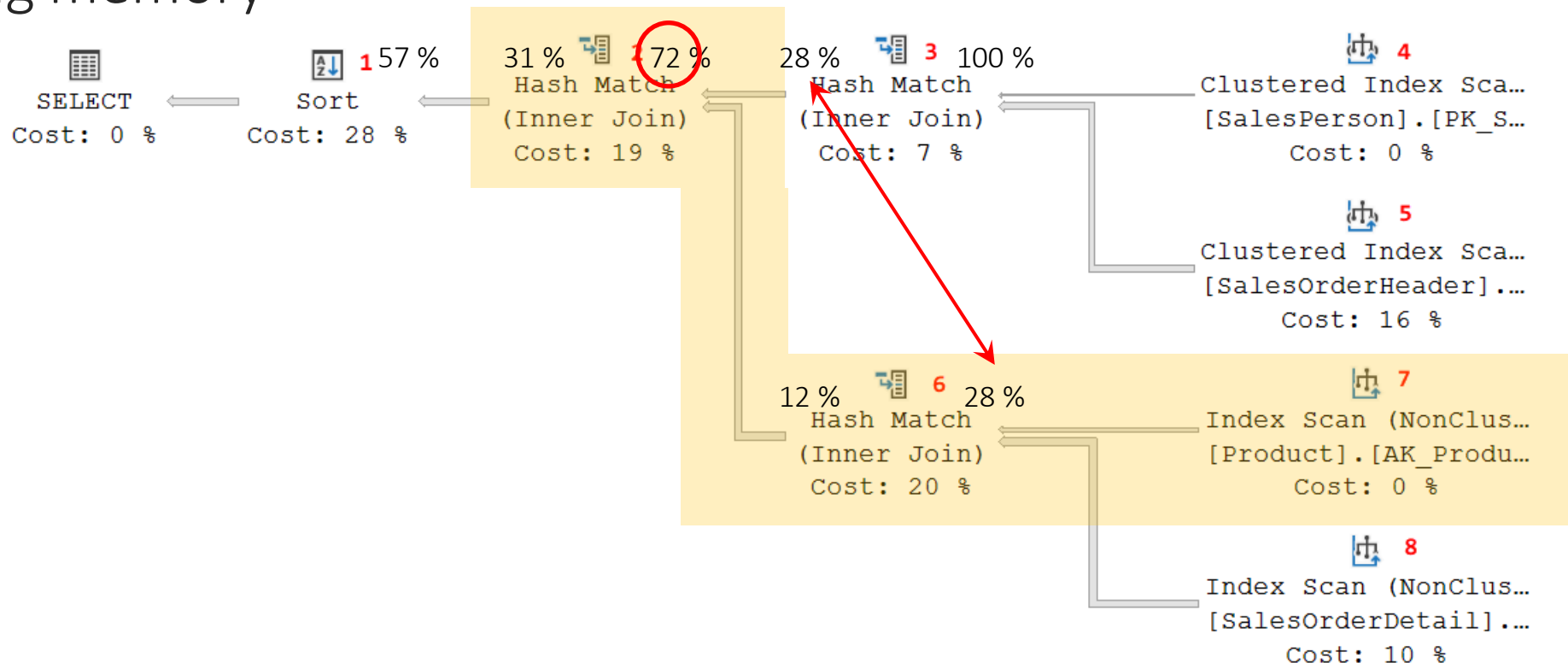
# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan
Reusing memory



SQLServerFast.com execution plan training, block 3, advanced level, chapter 3: Hash Match (advanced) - (c) Hugo Kornelis

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory



SELECT
Cost: 0 %

[1] Sort
Cost: 28 %

[2] Hash Match
(Inner Join)
Cost: 19 %

[3] Hash Match
(Inner Join)
Cost: 7 %    100 %

[4] Clustered Index Sca…
[SalesPerson].[PK_S…
Cost: 0 %

[5] Clustered Index Sca…
[SalesOrderHeader].…
Cost: 16 %

[6] Hash Match
(Inner Join)
Cost: 20 %

[7] Index Scan (NonClus…
[Product].[AK_Produ…
Cost: 0 %

[8] Index Scan (NonClus…
[SalesOrderDetail].…
Cost: 10 %

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan
Reusing memory

# Hash Match

Multiple *memory using* operators in a single execution plan

Reusing memory

Memory Fractions Input

Memory Fractions Output

# Hash Match

Multiple *memory using* operators in a single execution plan

    Reusing memory

    Hash teams

        Directly adjacent Hash match operators

        Using the same hash keys

        Team manager

            Mapping of values to buckets

            Mapping of buckets to partitions

            Memory management

            Decisions to spill

                Entire team spills at once

# Summary

Hash Match (advanced)
    Logic for all supported join types
    Hash spills
      Effect on order
    Memory fractions
    Hash teams

# Next chapters

Chapter 4: Adaptive Join (advanced)
    Logic for all supported join types
    Spills
    Considerations for choosing Adaptive Join