

Cost of accessing data (approximations)

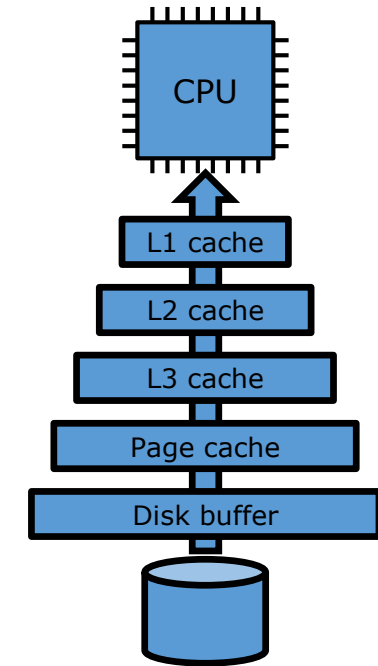
- Sequential reads
 - Option to maximize the effective read ratio
 - Depends on DB design
 - Enables pre-fetching

Cost = seek+rotation+n*transfer

- Random Access
 - Requires indexing structures
 - Ignores data locality

Cost_{single cylinder files} = seek+n*(rotation+transfer)

Cost_{multi-cylinder files} = n*(seek+rotation+transfer)



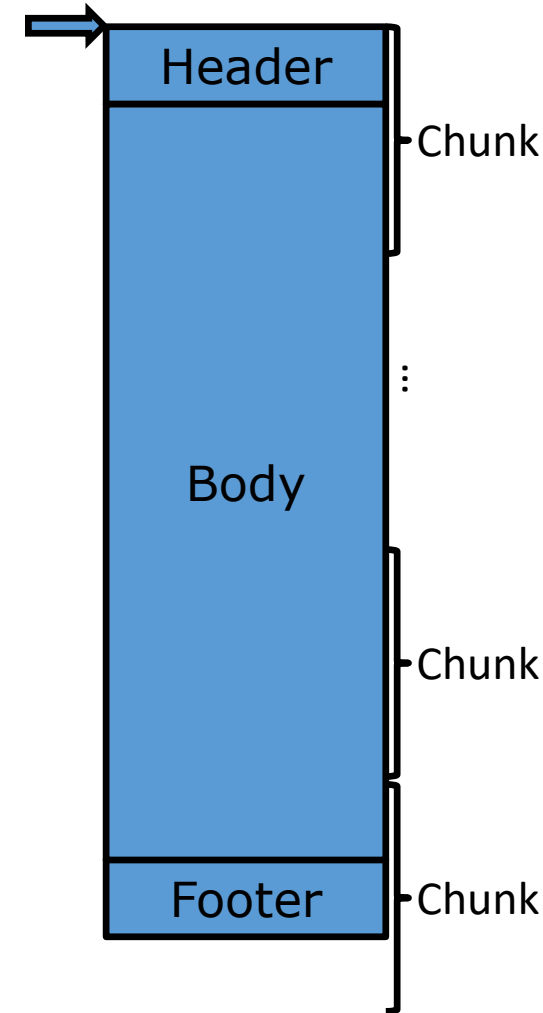
seek	~12ms
rotation	~3ms
transfer (8KB)	~0.03ms

General formulas for size estimation

$$\begin{aligned} \text{Size}(\text{Layout}) = & \text{Size}(\text{Header}_{\text{Layout}}) \\ & + \text{Size}(\text{Body}_{\text{Layout}}) \\ & + \text{Size}(\text{Footer}_{\text{Layout}}) \end{aligned}$$

$$\text{UsedChunks}(\text{Layout}) = \frac{\text{Size}(\text{Layout})}{\text{Size}(\text{chunk})}$$

$$\text{Seeks}(\text{Layout}) = \lceil \text{UsedChunks}(\text{Layout}) \rceil$$



Horizontal layout size estimation

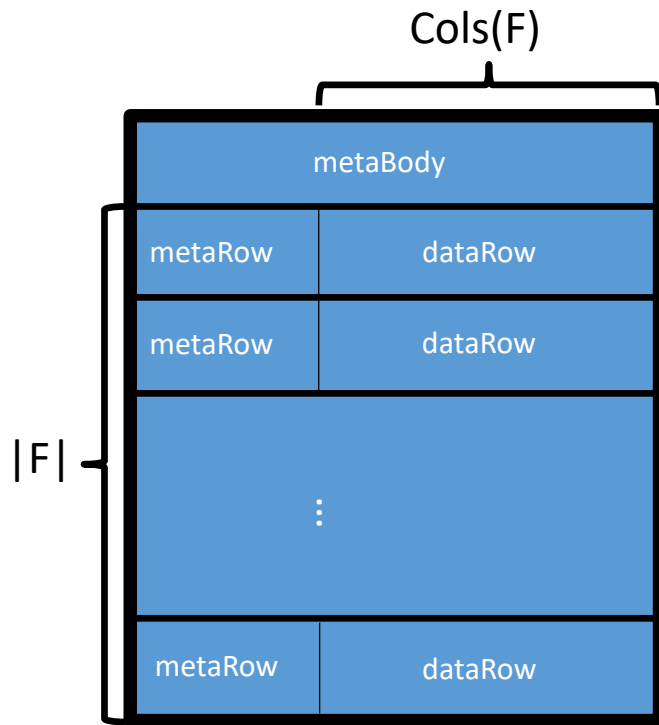


Table 1

A	B	C	D
101	201	301	401
102	202	302	402
103	203	303	403

Diagram illustrating the horizontal layout structure with specific values. The layout consists of a **metaBody** section at the top and a series of **metaRow** and **dataRow** pairs below it. The cardinality of the rows is denoted as $|F| = 3$, and the number of columns is denoted as $\text{Cols}(F) = 3$. A dashed arrow points from the **metaBody** section to the first cell of the first **dataRow**, which contains the value 101.

metaBody	
metaRow	101, 201, 301, 401
metaRow	102, 202, 302, 402
metaRow	103, 203, 303, 403

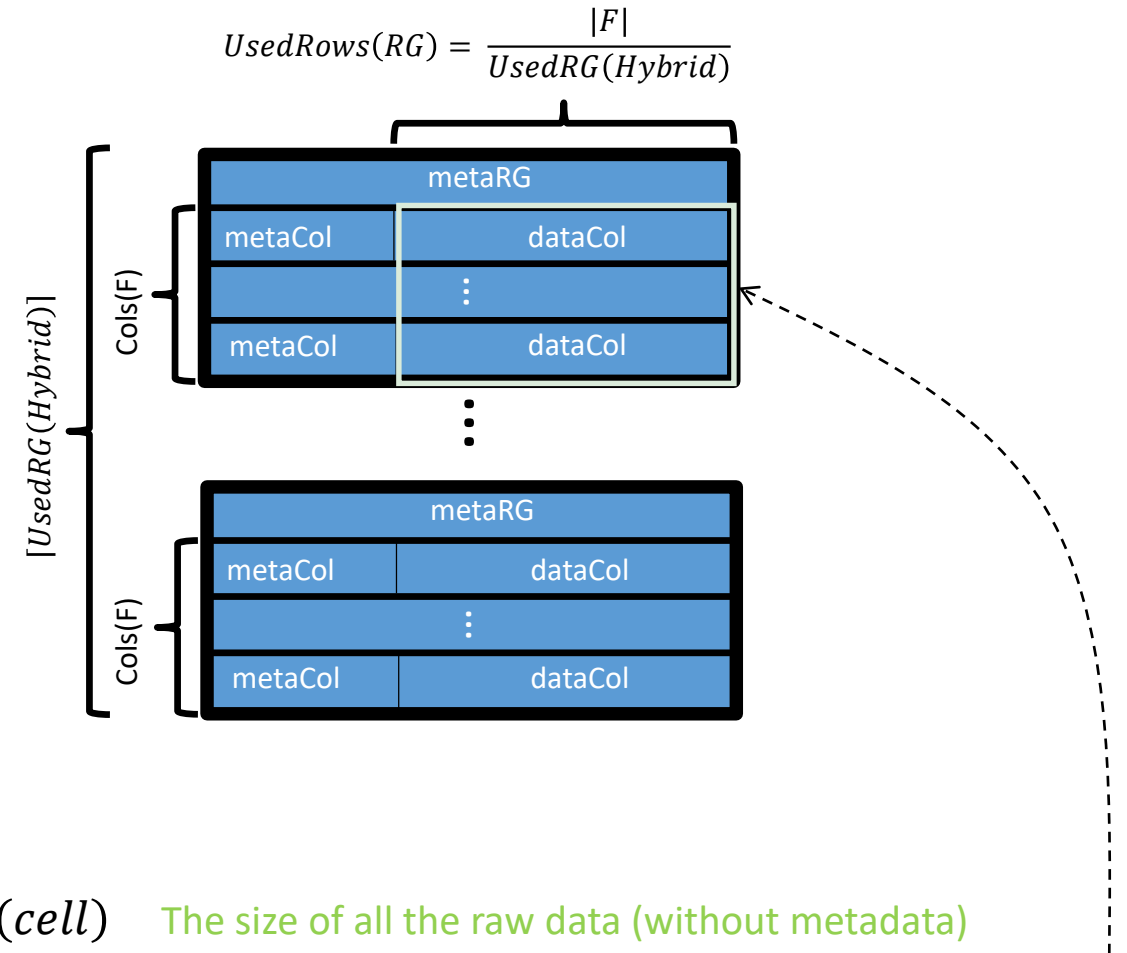
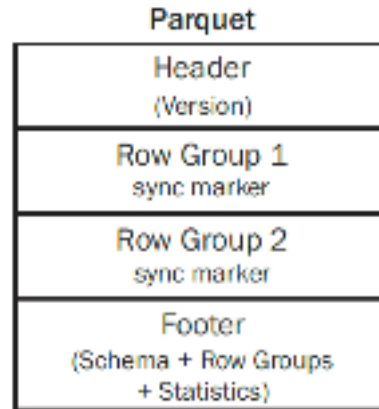
$|F| = 3$

What is the size of the body in Horizontal layouts?

$$\text{Size}(\text{dataRow}) = \text{Cols}(F) * \text{Size}(\text{cell})$$

$$\text{Size}(\text{Body}_{\text{Horizontal}}) = \text{Size}(\text{metaBody}) + |F| * (\text{Size}(\text{metaRow}) + \text{Size}(\text{dataRow}))$$

Hybrid layout size estimation



What is the size of the body in Vertical layouts?

$$UsedRG(Hybrid) = \frac{Cols(F) * |F| * Size(cell)}{(Size(RG) - Size(metaRG) - Cols(F) * Size(metaCol))}$$

The size of all the raw data (without metadata)

The size of an RG without metadata

$$Size(Body_{Hybrid}) = [UsedRG(Hybrid)] * (Size(metaRG) + Cols(F) * Size(metaCol)) + Cols(F) * |F| * Size(cell)$$

Multiply each RG with the metadata it stores

Plus the size of the raw data

Cost of projection in hybrid layouts

$$Size(projCols) = Proj(F) * UsedRows(RG) * Size(cell)$$

If UsedRows=dataCol
and Proj(F) = 2

$$Size(Project_{Hybrid}) = Size(Header_{Hybrid}) + Size(Footer_{Hybrid}) \\ + [UsedRG(Hybrid)] * (Size(metaRG) + proj(F) * Size(metaCol)) \\ + UsedRG(Hybrid) * Size(projCols)$$

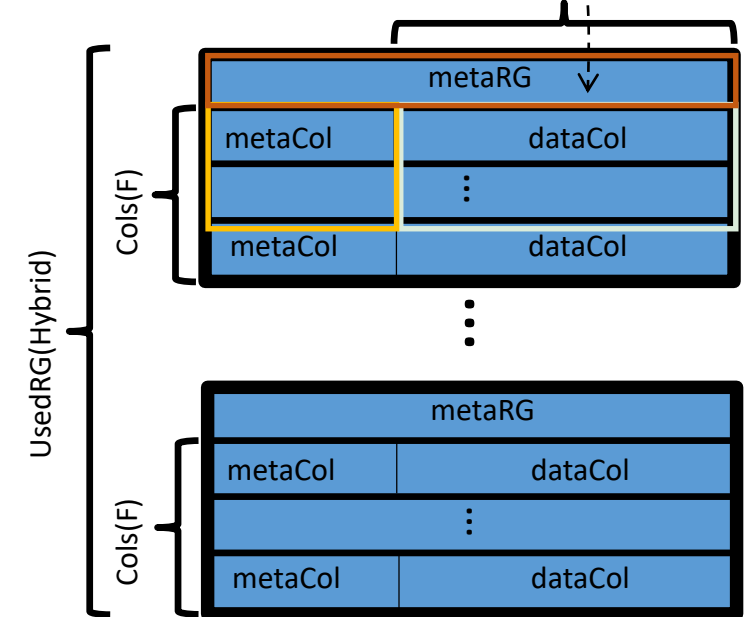
$$Cost(Project_{Hybrid}) = UsedChunks(Project_{Hybrid}) * W_{ReadTransfer} \\ + Seeks(Hybrid) * (1 - W_{ReadTransfer})$$

$$UsedRows(RG) = \frac{|F|}{UsedRG(Hybrid)}$$

What about selection?

We have to access the entire RG, even if only some data are actually requested by the user(e.g., there is a filter).

So, we need to compute how many RGs we need to access based on the rows that are defined by the predicate!



Selection in hybrid layouts

$$P(RGSelected) = 1 - (1 - SF)^{UsedRows(RG)}$$

$$Size(RowsSelected) = \left\lceil \frac{SF * |F|}{UsedRows(RG)} \right\rceil (Size(metaRG) + Cols(F) * Size(metaCol))$$

The size of the metadata depending on SF

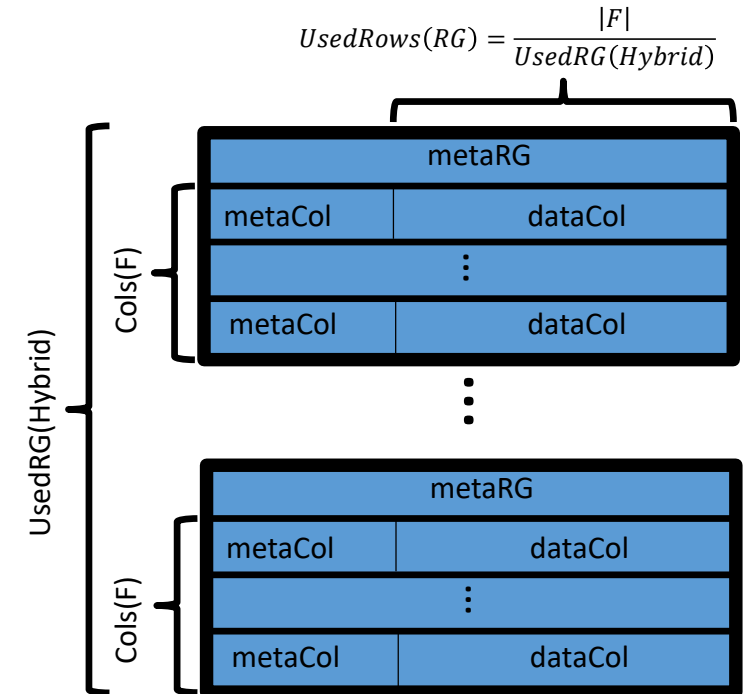
$$+ SF * |F| * Cols(F) * Size(cell)$$

The size of the raw data multiplied by SF

$$UsedRG(Select_{Hybrid}) = \begin{cases} \text{if unsorted: } P(RGSelected) * UsedRG(Hybrid) \\ \text{if sorted: } \frac{Size(RowsSelected)}{Size(RG)} \end{cases}$$

$$Size(Select_{Hybrid}) = Size(Header_{Hybrid}) + Size(Footer_{Hybrid}) + UsedRG(Select_{Hybrid}) * Size(RG)$$

$$Cost(Select_{Hybrid}) = UsedChunks(Select_{Hybrid}) * W_{ReadTransfer} + Seeks(Select_{Hybrid}) * (1 - W_{ReadTransfer})$$



Client caching

Cash miss

1. The client sends a READ command to the coordinator
2. The coordinator requests chunkservers to send the chunks to the client
 - Ranked according to the closeness in the network
3. The list of locations is cached in the client
 - Not a complete view of all chunks

Cash hit

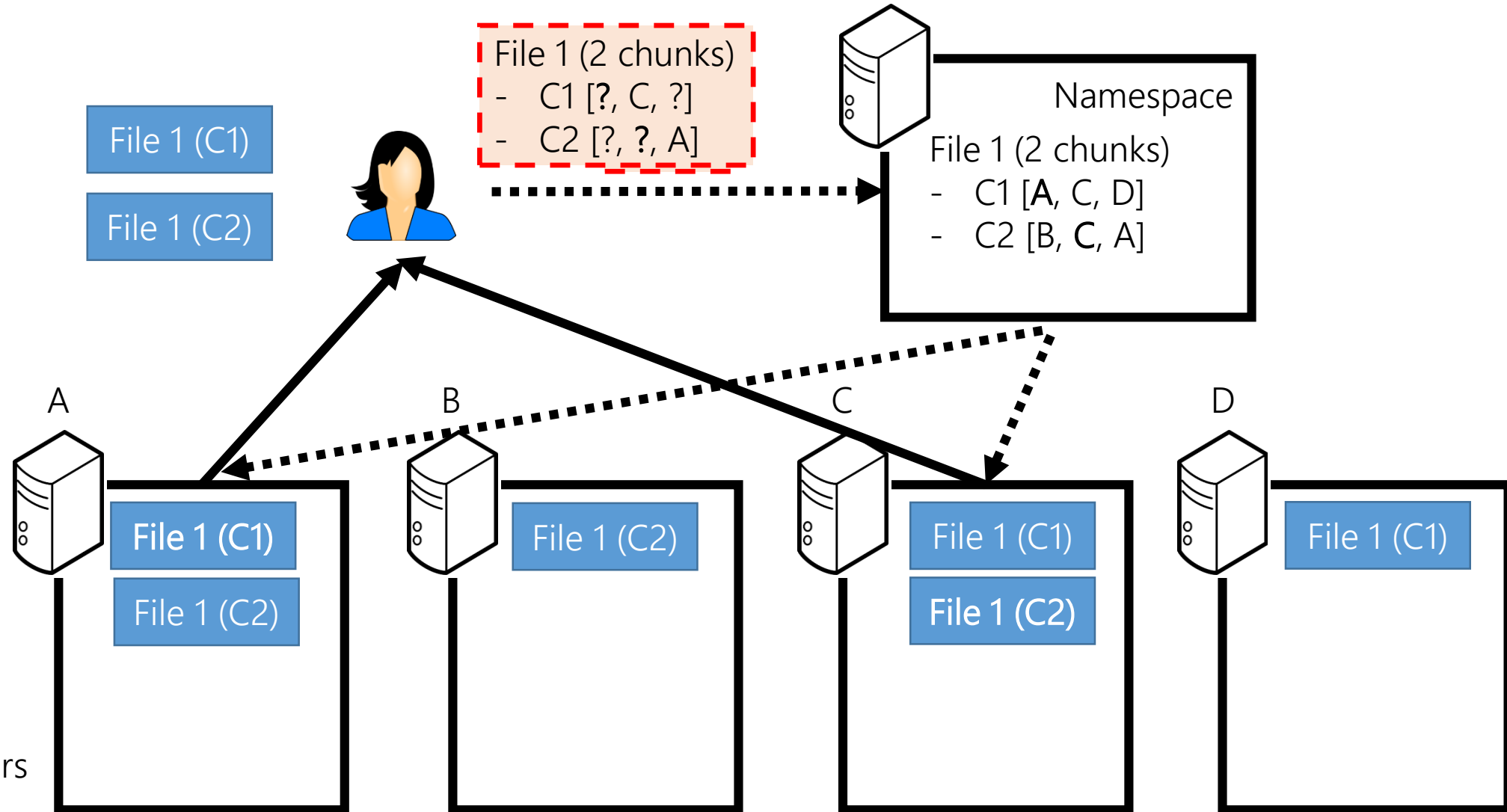
1. The client reads the cache and requests the chunkservers to send the chunks

Avoid coordinator bottleneck
+
One communication step is saved

Example of client cache miss

Coordinator

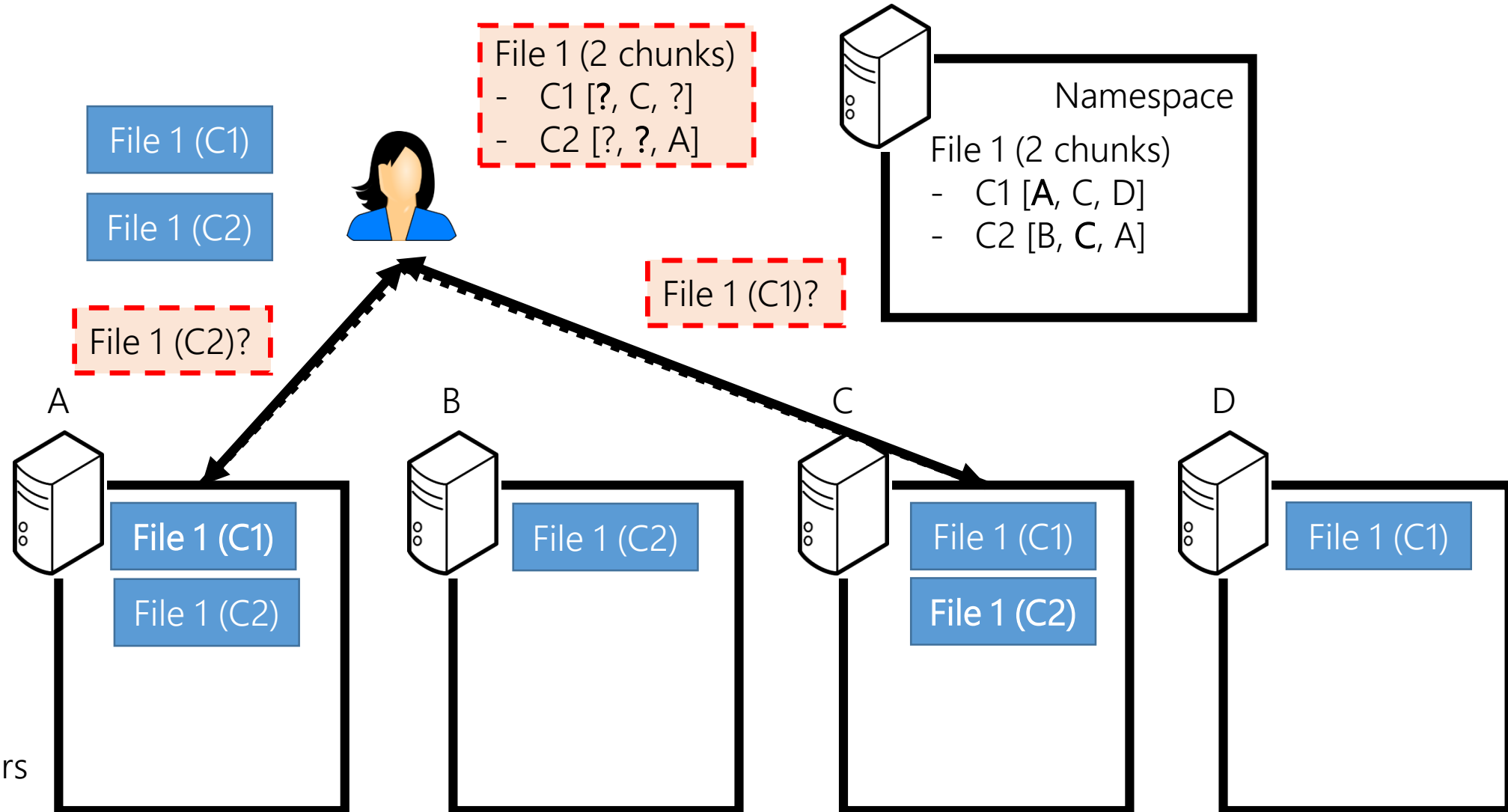
—————▶ Data flow
.....▶ Control flow



ChunkServers

Example of client cache hit

—————▶ Data flow
.....▶ Control flow



ChunkServers