

Dynamic Indexing with Logarithmic Merging

Alfan F. Wicaksono

Fakultas Ilmu Komputer, Universitas Indonesia

Dynamic Indexing using Logarithmic Merging

We can do better than $O(T^2/n)$ by introducing $\log_2(T/n)$ indexes I_0, I_1, I_2, \dots of size $2^0 \times n, 2^1 \times n, 2^2 \times n, \dots$

Postings percolate up this sequence of indexes and are processed only once on each level.

LMERGEADDTOKEN(indexes, Z_0 , token):

$Z_0 \leftarrow \text{MERGE}(Z_0, \{\text{token}\})$

if $|Z_0| = n$ **then**

for $i \leftarrow 0$ **to** ∞ **do**

if $I_i \in \text{indexes}$ **then**

$Z_{i+1} \leftarrow \text{MERGE}(I_i, Z_i)$ (Z_{i+1} is a temporary index on disk)

$\text{indexes} \leftarrow \text{indexes} - \{I_i\}$

else

$I_i \leftarrow Z_i$ (Z_i becomes the permanent index I_i)

$\text{indexes} \leftarrow \text{indexes} \cup \{I_i\}$

BREAK

$Z_0 \leftarrow \emptyset$

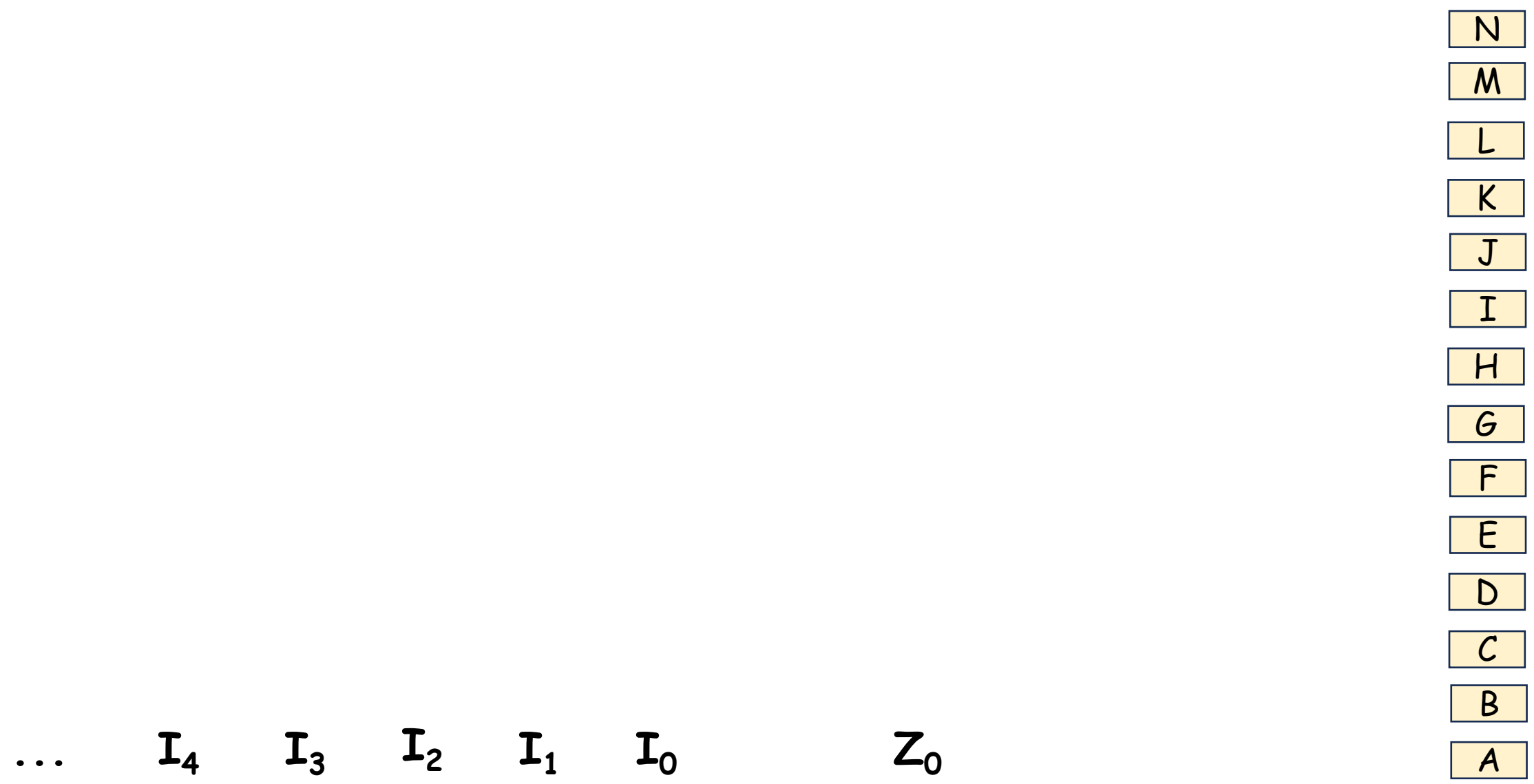
LOGARITHMICMERGE():

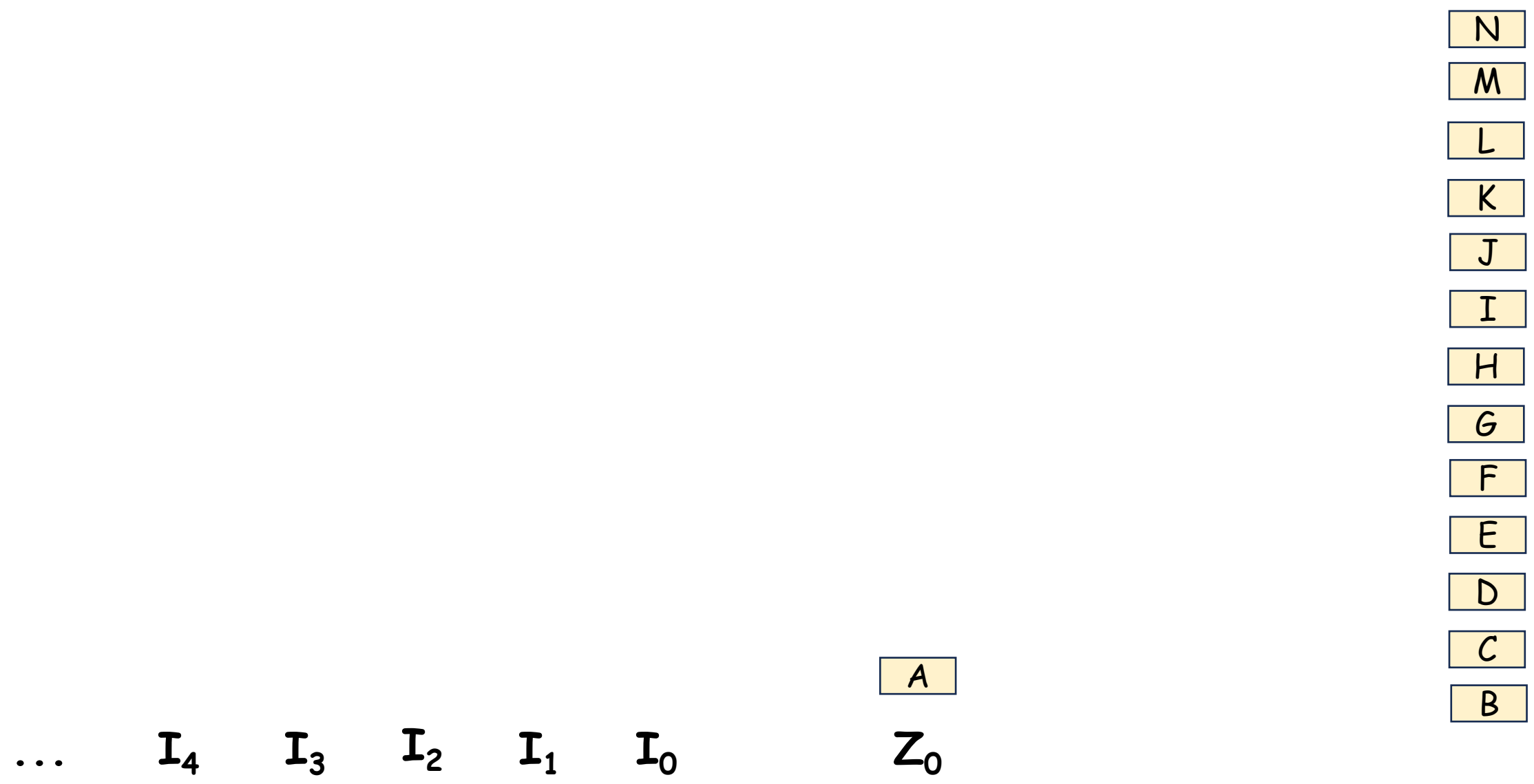
$Z_0 \leftarrow \emptyset$ (Z_0 is the in-memory index.)

$\text{indexes} \leftarrow \emptyset$

while true **do**

$\text{LMERGEADDTOKEN}(\text{indexes}, Z_0, \text{GETNEXTTOKEN}())$





- N
- M
- L
- K
- J
- I
- H
- G
- F
- E
- D
- C

B

A

...

I_4

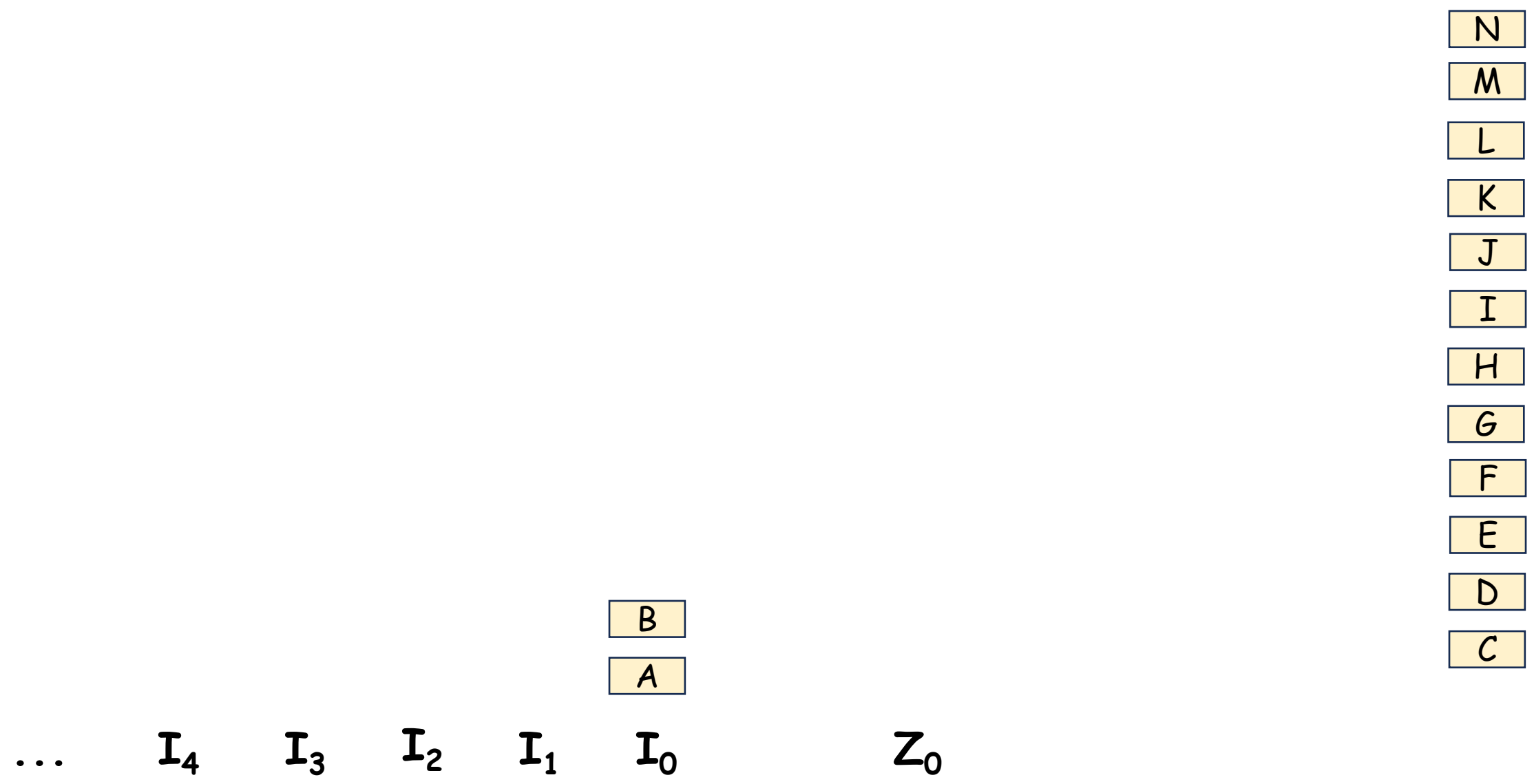
I_3

I_2

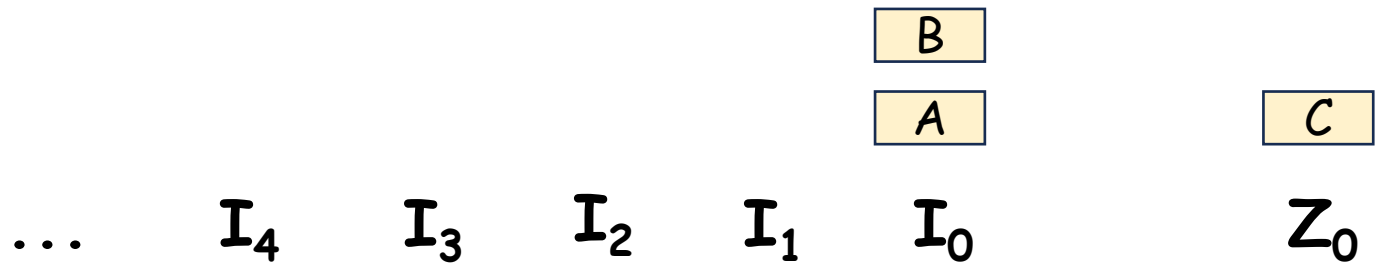
I_1

I_0

Z_0



- N
- M
- L
- K
- J
- I
- H
- G
- F
- E
- D



- N
- M
- L
- K
- J
- I
- H
- G
- F
- E

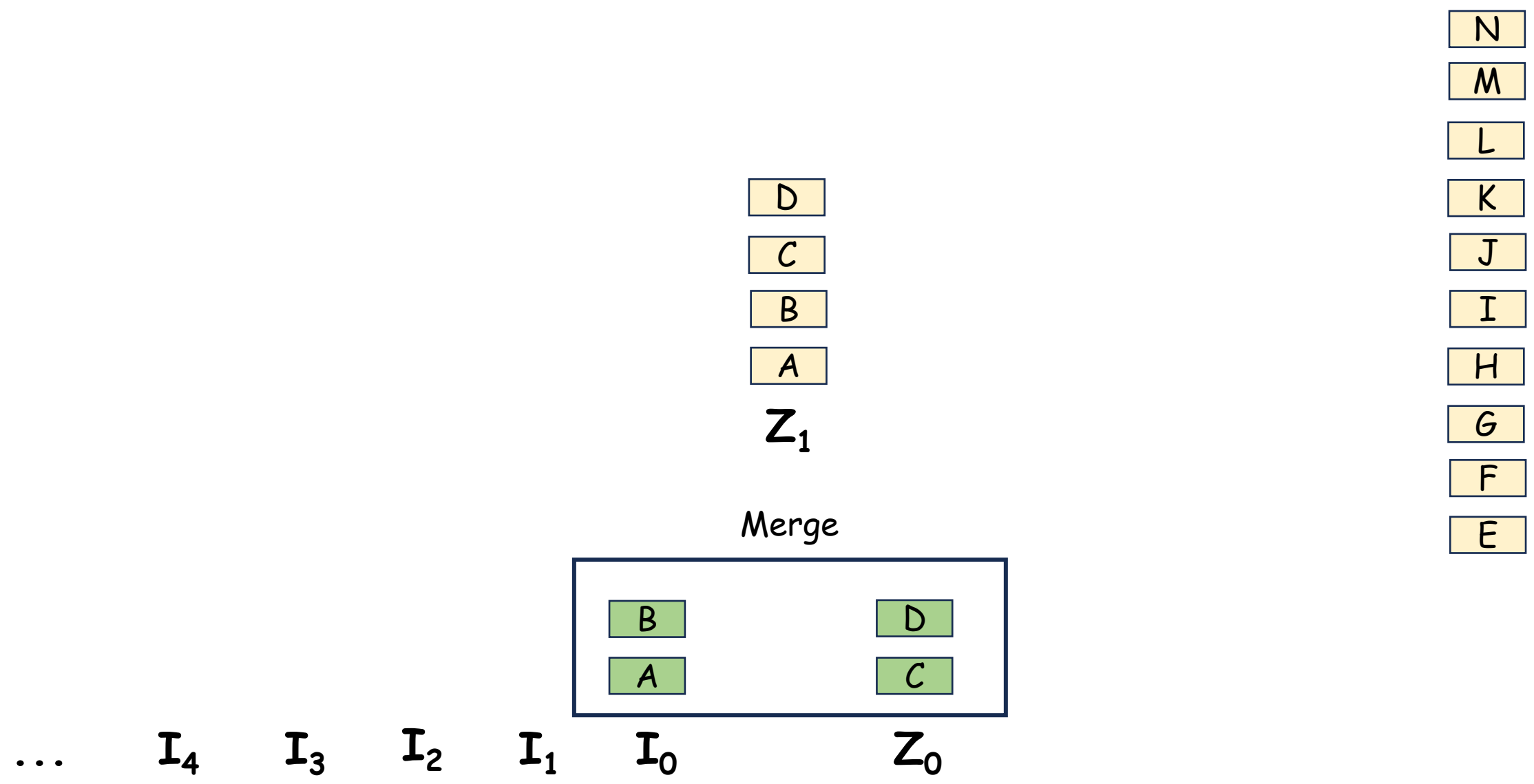
B

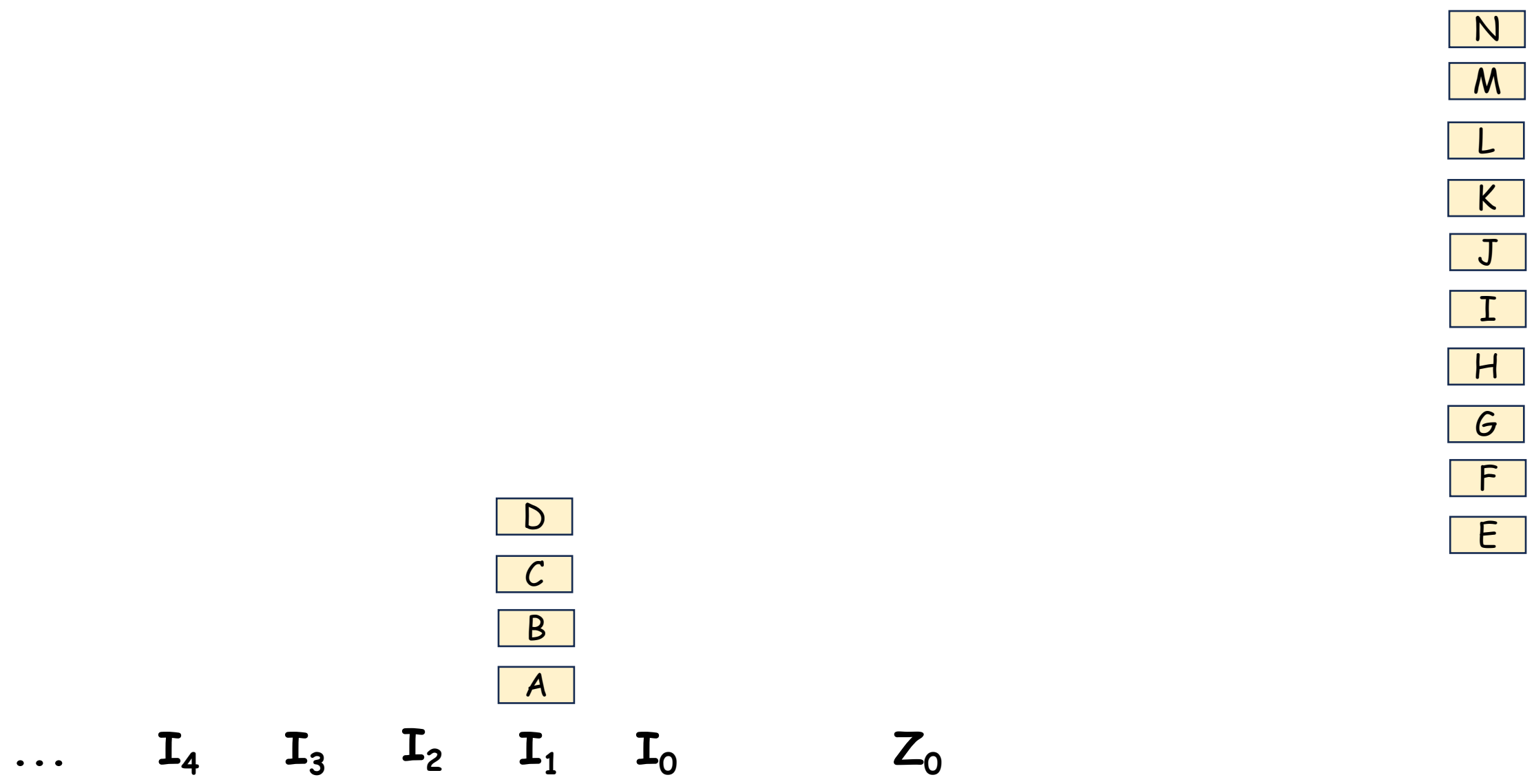
A

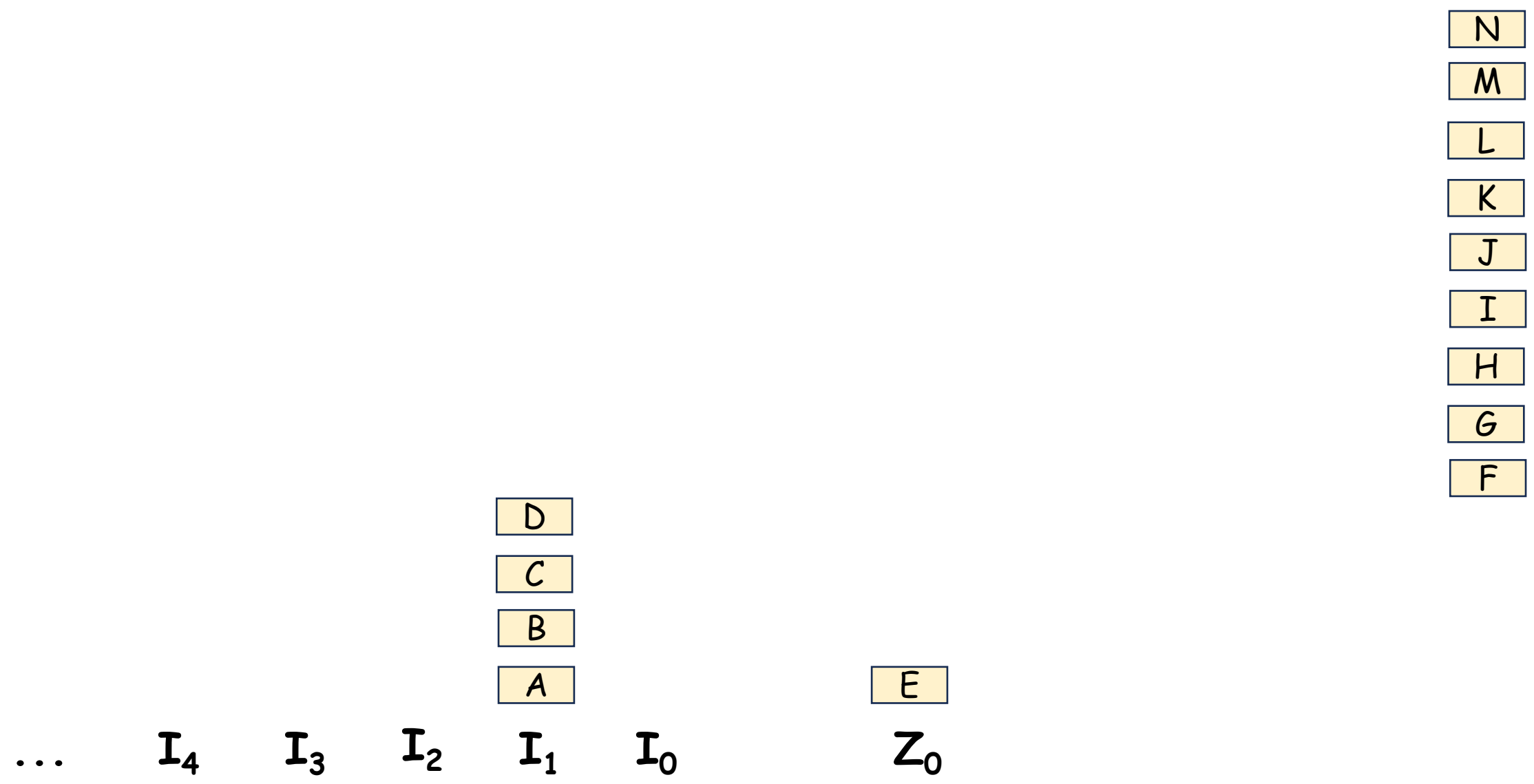
D

C

... I_4 I_3 I_2 I_1 I_0 Z_0







N
M
L
K
J
I
H
G

D
C
B
A

F
E

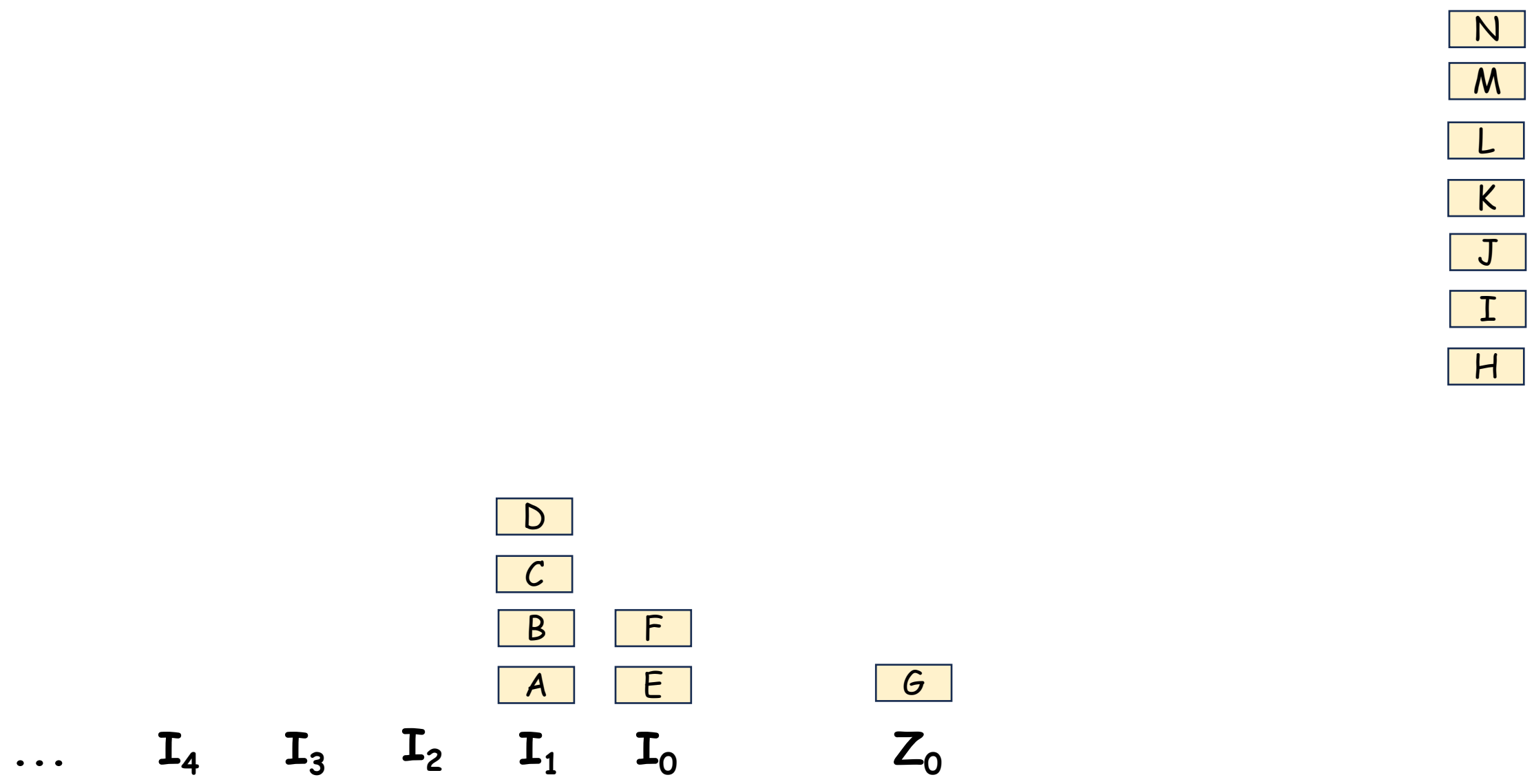
... I_4 I_3 I_2 I_1 I_0 Z_0

- N
- M
- L
- K
- J
- I
- H
- G

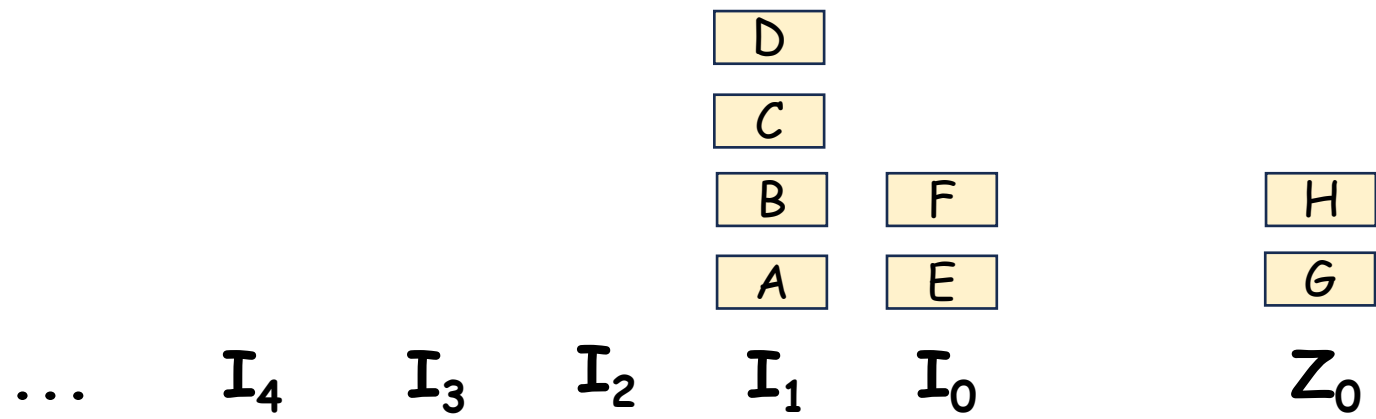
- D
- C
- B
- A

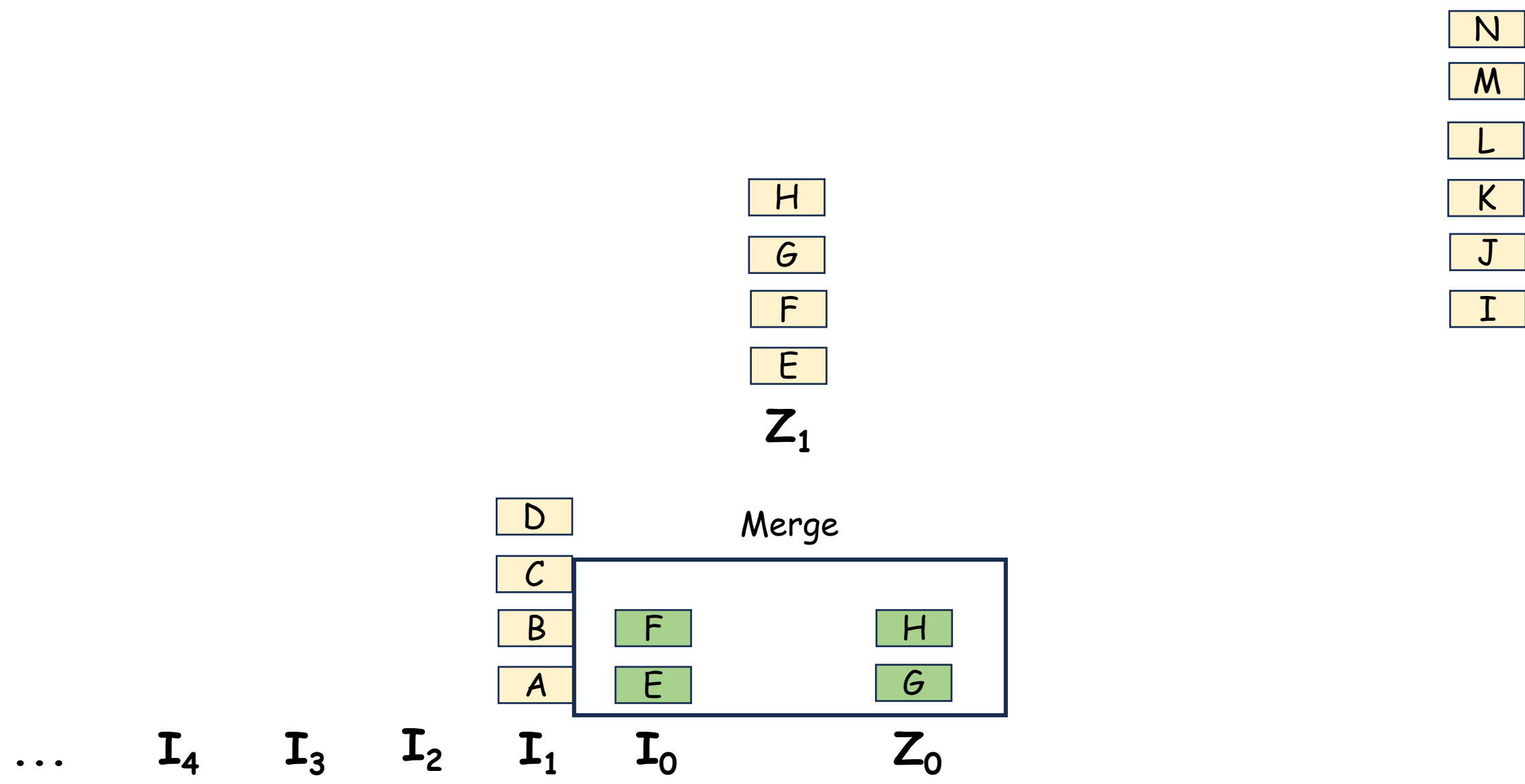
- F
- E

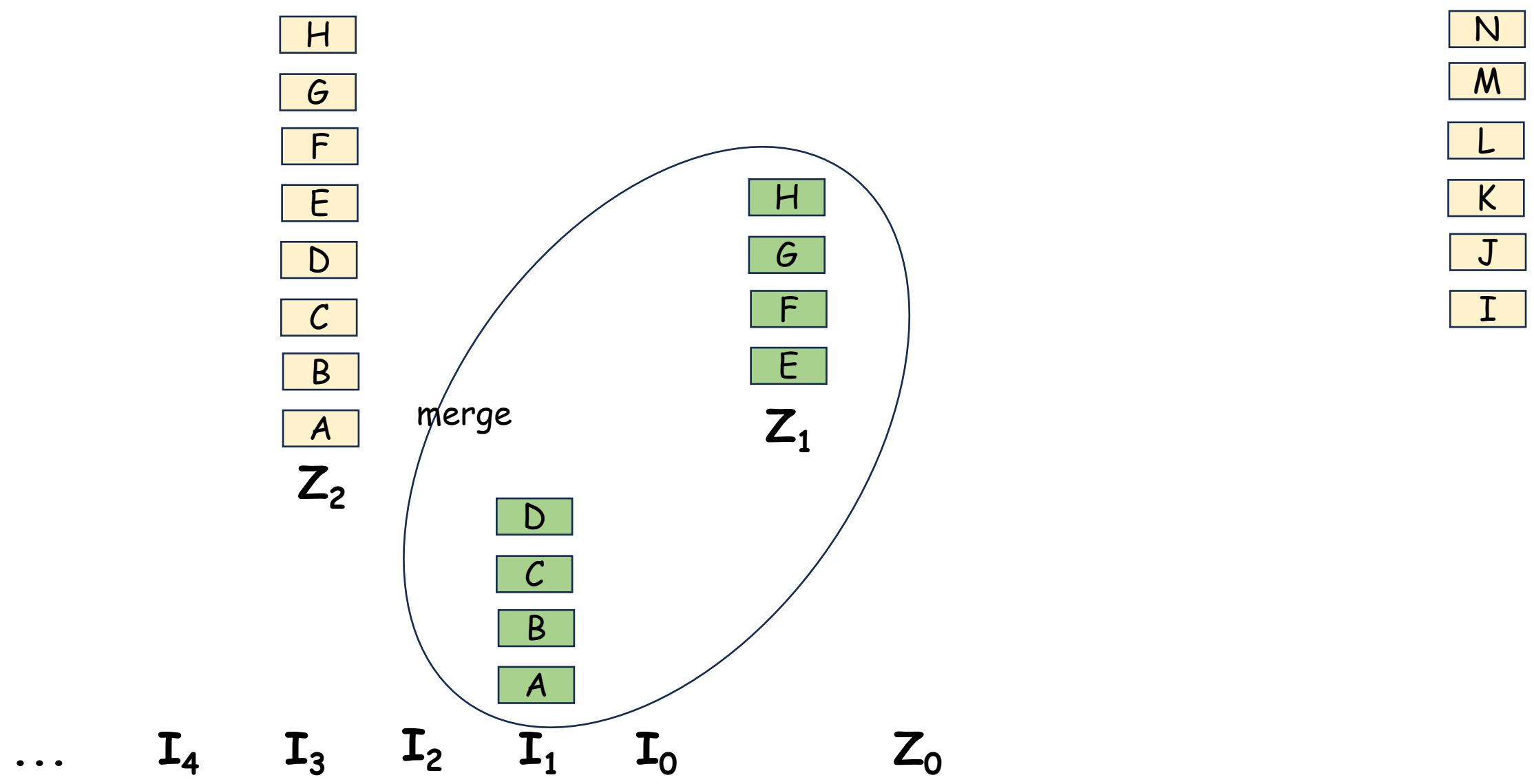
... I_4 I_3 I_2 I_1 I_0 Z_0

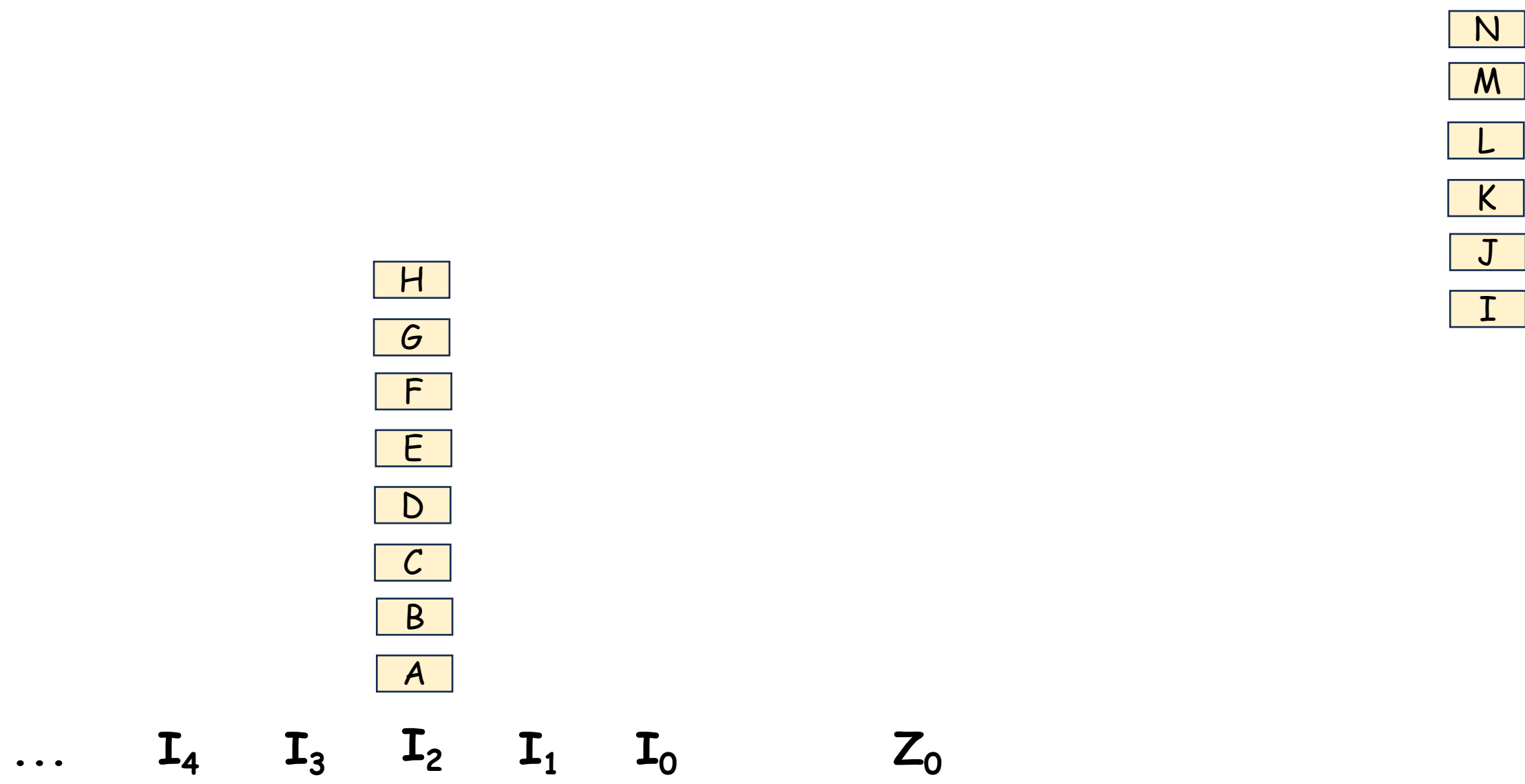


N
M
L
K
J
I









N
M
L
K
J

H
G
F
E
D
C
B
A

I

... I_4 I_3 I_2 I_1 I_0 Z_0

N
M
L
K

H
G
F
E
D
C
B
A

J
I

... I_4 I_3 I_2 I_1 I_0 Z_0

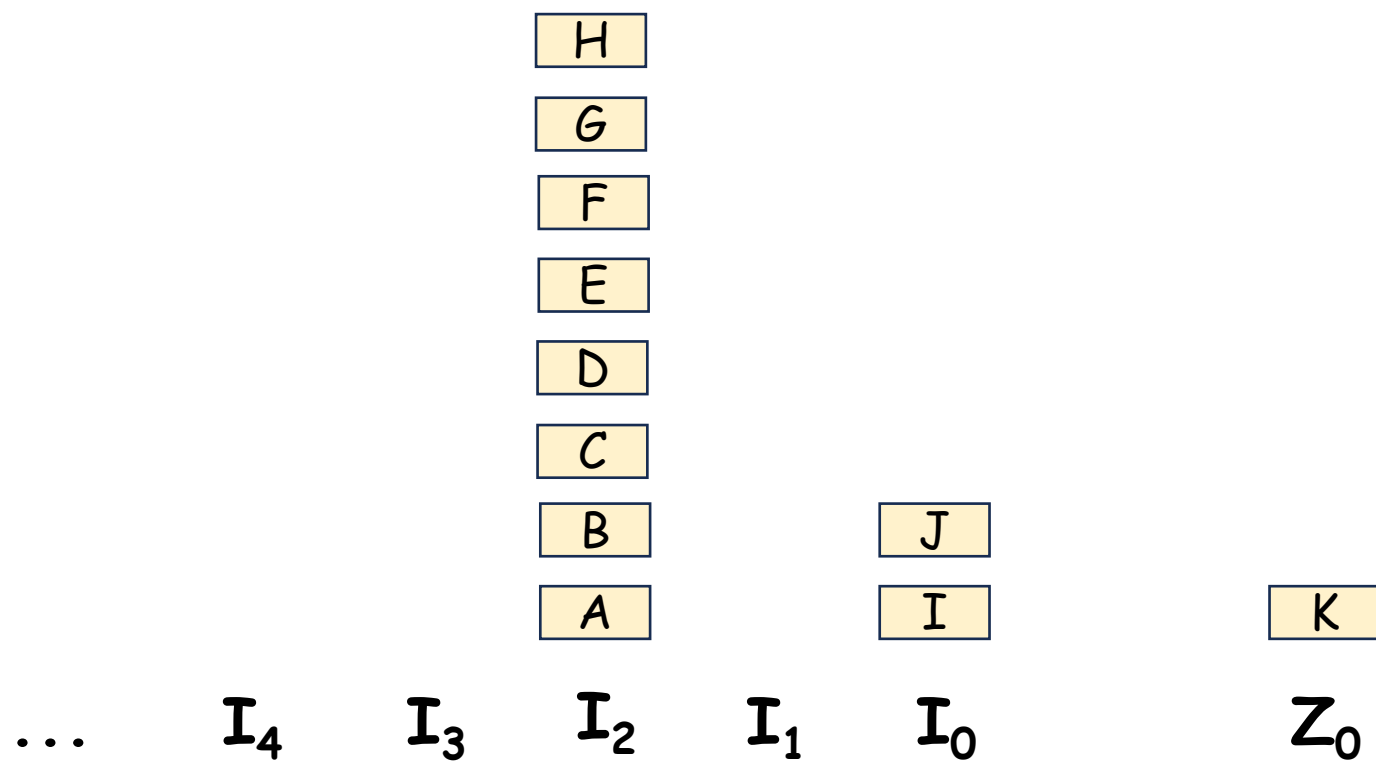
N
M
L
K

H
G
F
E
D
C
B
A

J
I

... I_4 I_3 I_2 I_1 I_0 Z_0

N
M
L



N

M

H

G

F

E

D

C

B

A

J

I

L

K

...

I_4

I_3

I_2

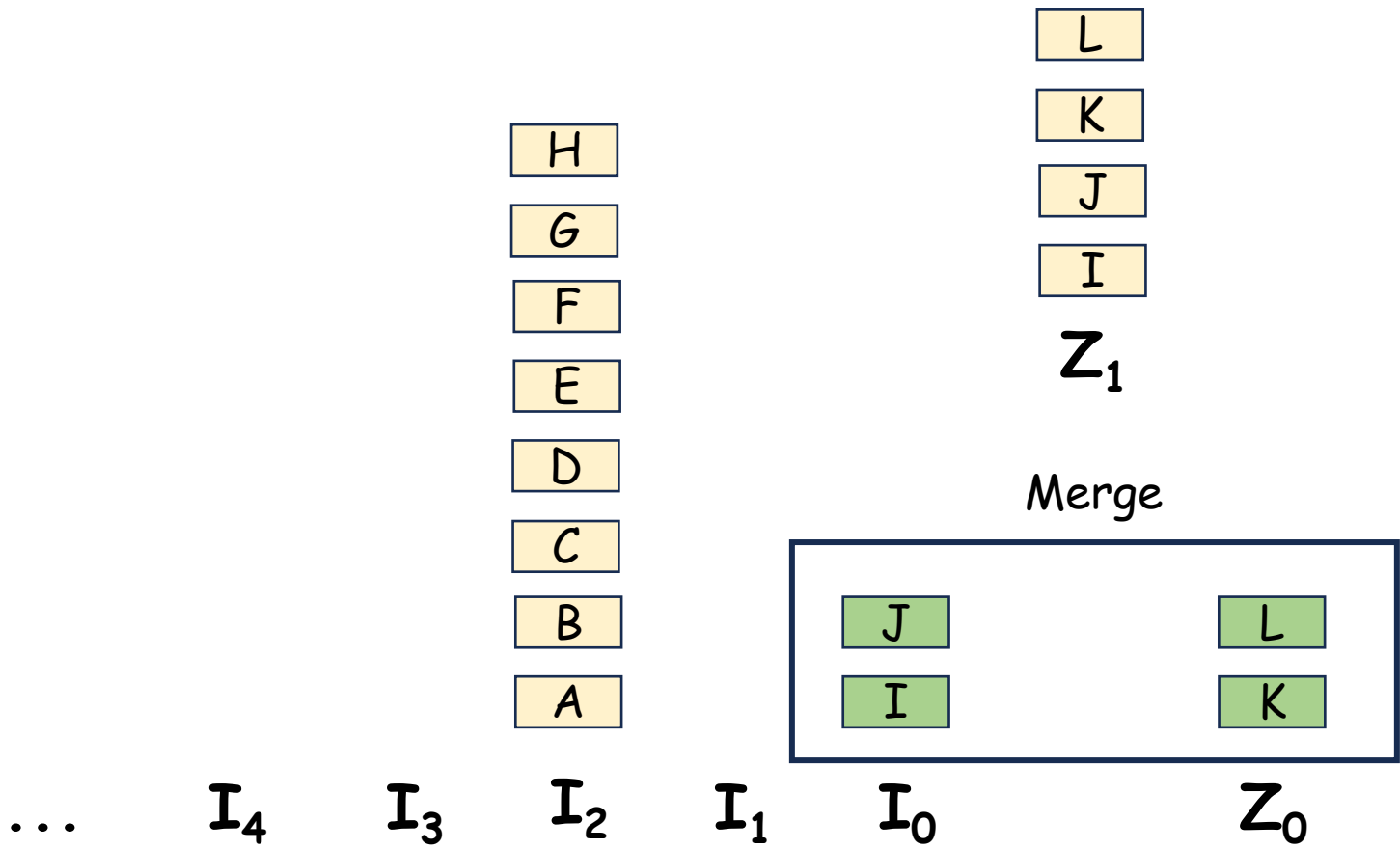
I_1

I_0

Z_0

N

M



N

M

H

G

F

E

D

C

B

A

L

K

J

I

...

I_4

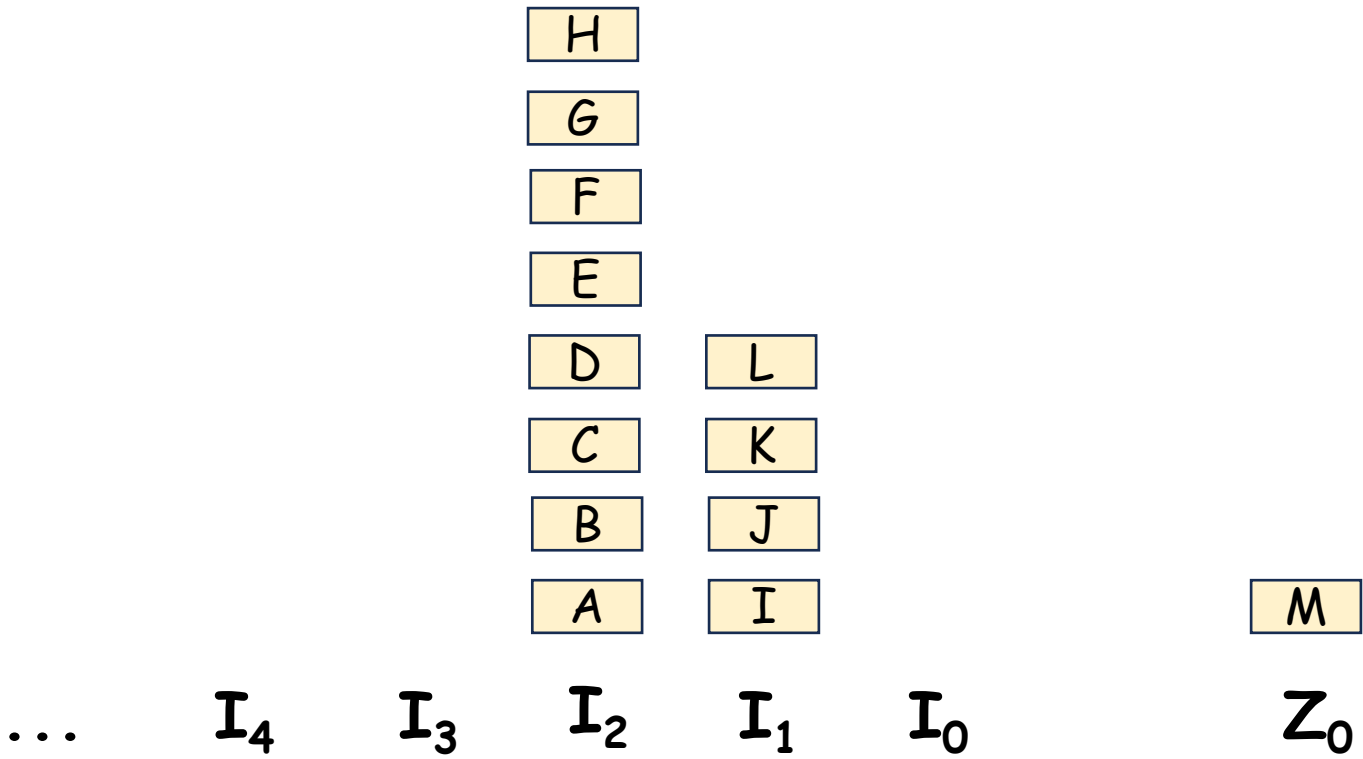
I_3

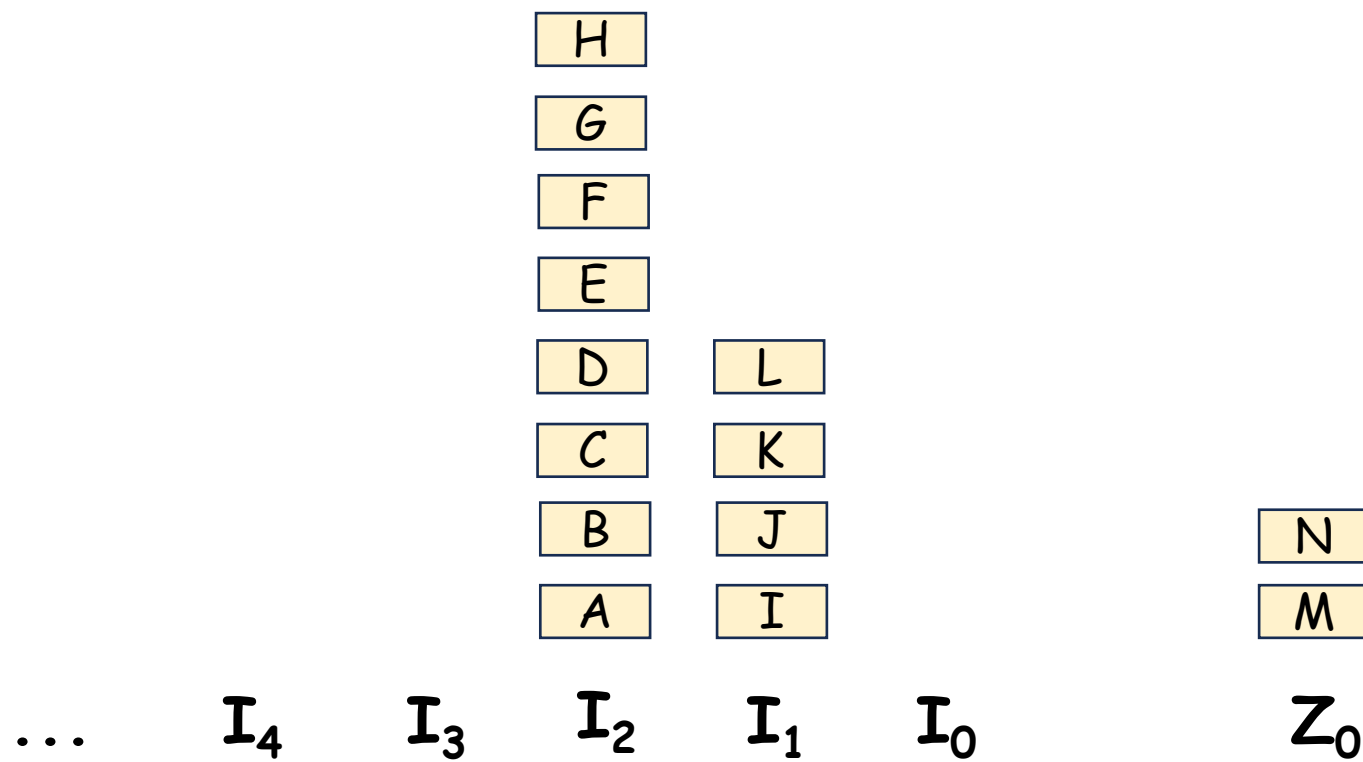
I_2

I_1

I_0

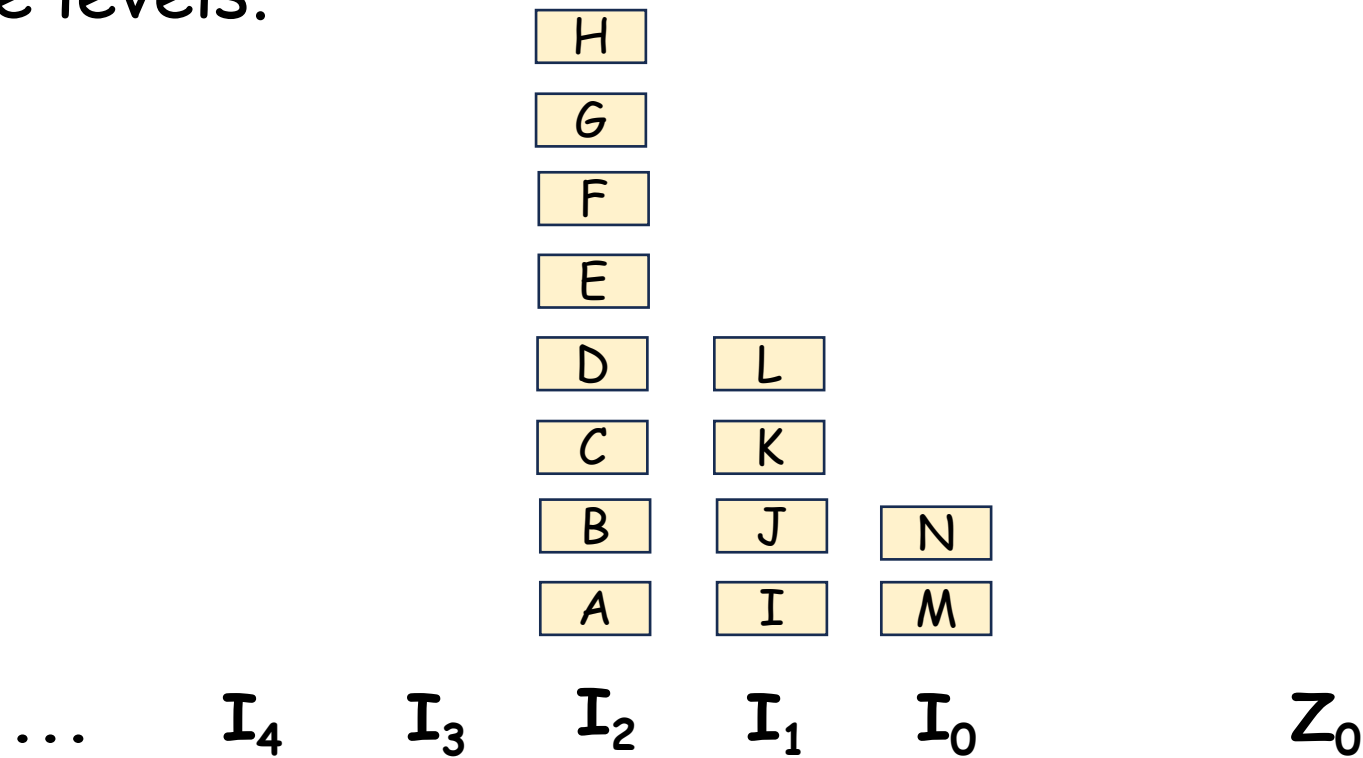
Z_0





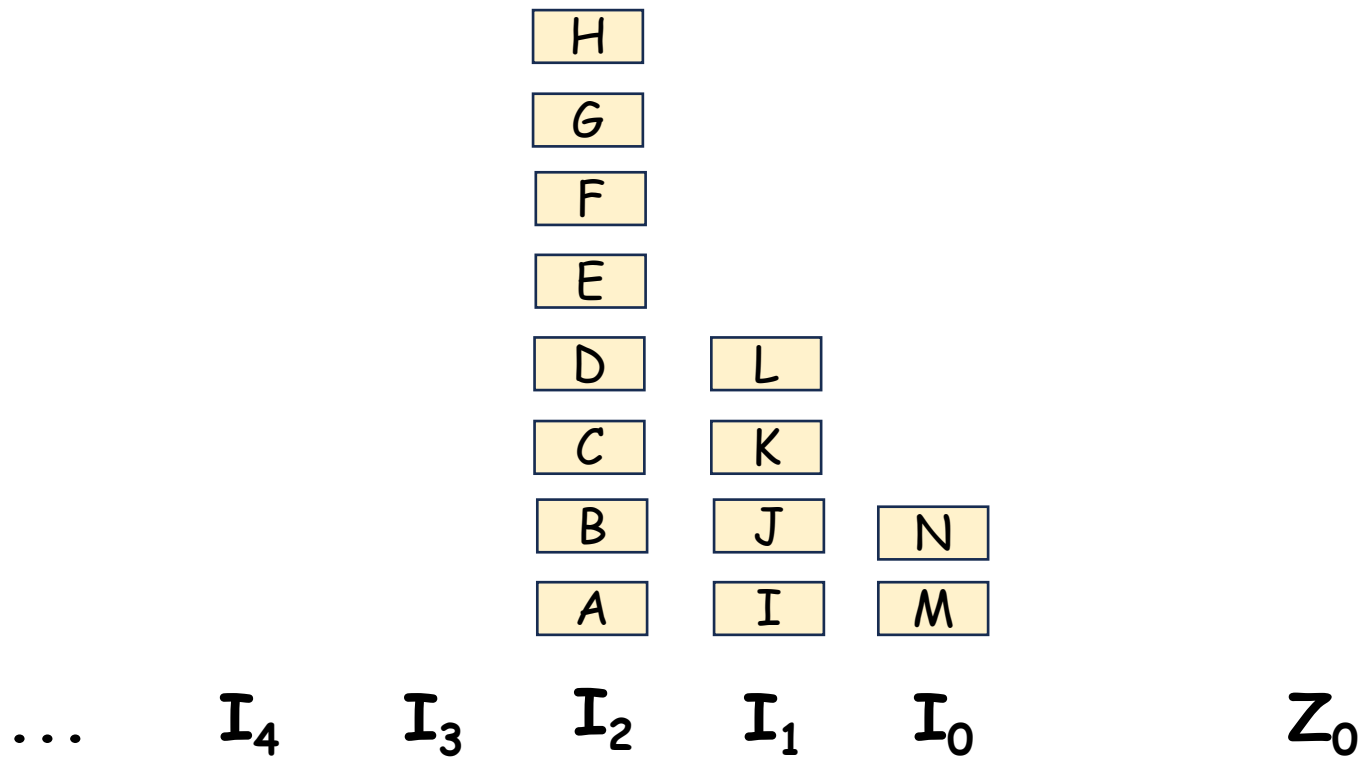
There are $\log_2(T/n)$ indexes.

Overall index construction time is $O(T \log_2(T/n))$ because each posting is processed only once on each of $\log_2(T/n)$ the levels.



Dan seterusnya ...

However, We trade this efficiency gain for a slow down of query processing. We now need to merge results from $\log_2(T/n)$ indexes as opposed to **just two**.



Dan seterusnya ...