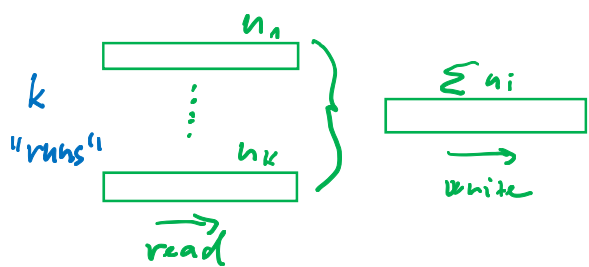


# Multway mergesort

(better use of cache)



needs  $\lceil \log_k N \rceil = \lceil \log N / \log k \rceil$  passes  
 minimum of  $k$  items by heap  $\Rightarrow \Theta(\log k)$  per item  
 single pass  $\Theta(N \log k)$  time

$\Rightarrow$   $k$ -way mergesort in  $\Theta(N \log k \cdot \log N / \log k) = \Theta(N \log N)$  time  
 ↑ independent on  $k$

## I/O complexity:

single merge:  $\Theta(\sum n_i / B + 1)$  ↑ runs are consecutive

single pass:  $\Theta(N/B + 1)$  ↑ successive merges on consecutive runs

in total:  $\Theta(N/B \log N / \log k + 1)$  ↑ for  $N < B$  when  $O(1)$  suffice otherwise  $O(N/B + 1) = O(N/B)$

## How large cache we need?

$k+1$  blocks for single merge, heap  $\leq k-1$  blocks  $\Rightarrow M \geq 2Bk$  suffice

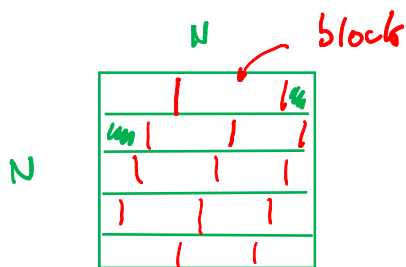
- In I/O model + cache-aware model, set  $k := \lfloor M/2B \rfloor$

$\Rightarrow \Theta(N/B \frac{\log N}{\log M/B} + 1)$  block transfers (optimal)

- In cache-oblivious model this cannot be done. (funnel sort)

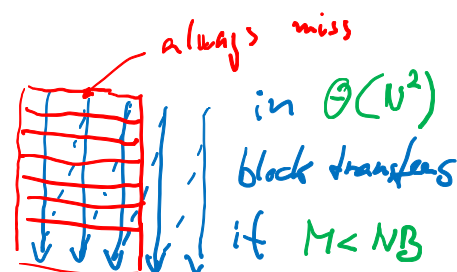
## Matrix transposition

(similarly matrix multiplication, ...)



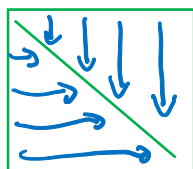
traversing the matrix

in  $O(N^2/B + 1)$  block transfers  
 (scanning array)



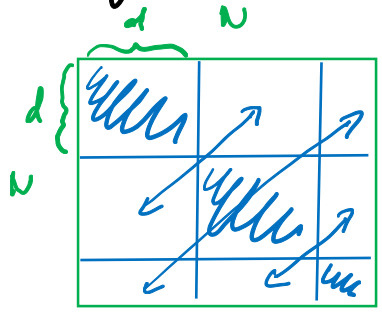
row-major order

$\Rightarrow$  naive transposition has  $\Theta(N^2)$  block transfers if  $M = o(NB)$



(same problem occurs in a constant fraction of columns)

Tiling (for cache-aware alg.)



tiles  $d \times d$  (or smaller on borders) in  $N \times N$  matrix

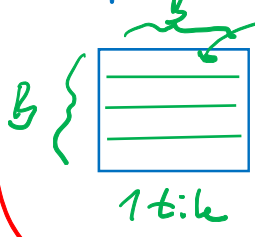
#tiles =  $\lceil N/d \rceil^2 \leq (N/d+1)^2 = O(N^2/d^2+1)$

Alg: transpose diag. tiles, transpose+swap pairs of non-diagonal tiles

Idea: if 2 tiles fit into the cache, subproblems can be solved recursively.

$\Rightarrow$  set  $d := B \Rightarrow O(N^2/B^2+1)$  tiles

1) perfect alignment ( $B$  divides  $N$ )



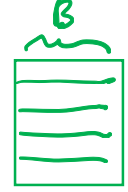
each tile needs  $B$  blocks

$O(N^2/B^2)$  tiles  $\Rightarrow O(N^2/B)$  block transfers

How large cache we need?

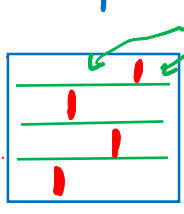
$M \geq 2B^2$

"Tall cache assumption":

$M = \Omega(B^2) \geq \epsilon B$  

if const.  $\epsilon < 2$ , use tiles const. times smaller  $\Rightarrow$  same asymptotics

2) not perfect alignment ( $B$  does not divide  $N$ )



each tile needs  $2B$  block  $\Rightarrow O(B)$  transfers

$O(N^2/B^2+1)$  tiles  $\Rightarrow O(N^2/B+1)$  block transfers

How large cache we need?

$M \geq 4B^2$

Alg. needs to know  $B$  (the tile size)  $\Rightarrow$  cache-aware.

Runs in  $O(N^2)$  time and  $O(N^2/B+1)$  block transfers.

optimal as each non-diagonal element has to be transferred

# Cache oblivious matrix transposition

cache-aware strategy: decompose into subproblems that fit into cache

vs

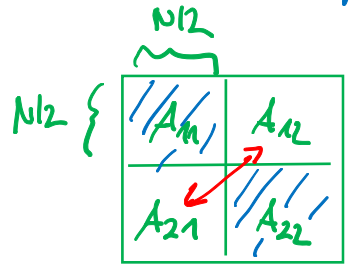
cache-oblivious strategy: alg. does not know B.M but cache controller does!

- recursively subdivide to smaller and smaller problem sizes
- on some level of recursion, the cache controller employs

the **same** replacement strategy as for cache aware alg. =>

**same** I/O complexity (up to multiplicative constants)

Divide and conquer alg. 1) assume  $N=2^k$  for some integer  $k$



"naive" approach: transpose each submatrix

$$T_N \rightarrow 4T_{N/2} + S_{N/2}$$

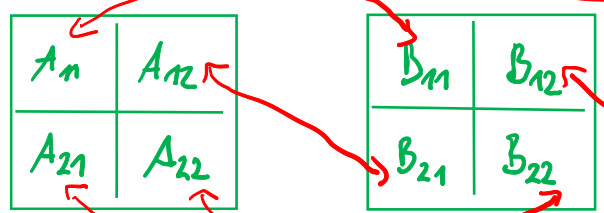
+ swap  $A_{21}, A_{12}$

problem size

by (note) scanning ✓

=>  $O(N^2 \log N)$  (spoiled) time

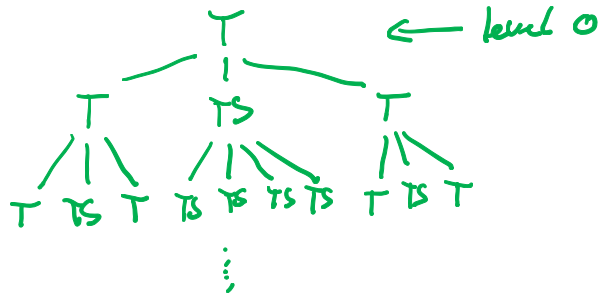
"transpose + swap" operation  $TS(A, B)$ :



$$T_N \rightarrow 2T_{N/2} + TS_{N/2}$$

$$TS_N \rightarrow 4TS_{N/2}$$

tree of recursion:



height  $k = \log N$  (last level)

# leaves  $\leq 4^{\log N} = N^2$

# internal nodes  $<$  # leaves  $\leq N^2$

leaves:  $T(a_{ii})$  do nothing,  $TS(a_{ij}, a_{ji})$  swap in  $O(1)$  time

internal nodes: redistribute work (compute indices of submatrices)

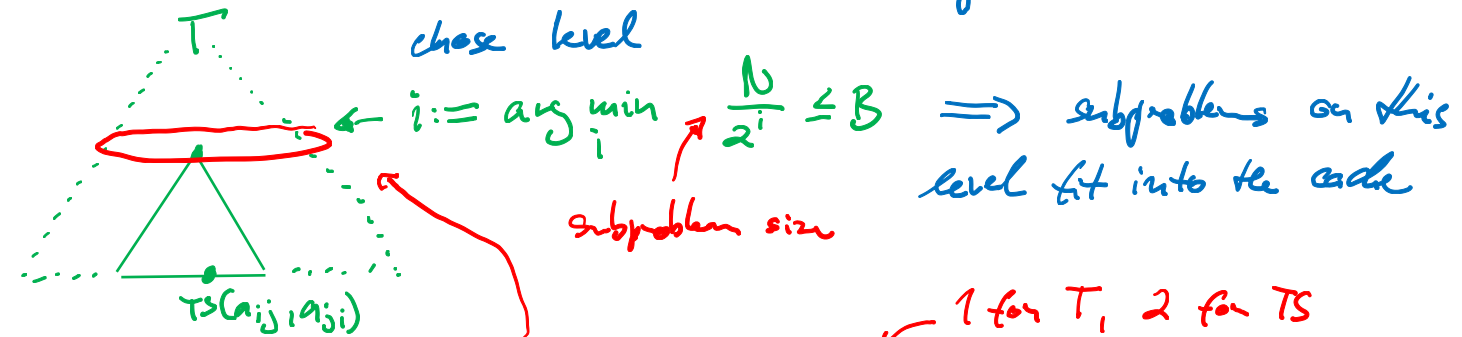
=>  $O(1)$  time per node =>  $O(N^2)$  total time (as needed)

I/O analysis

Is there a "good" replacement strategy for this alg.?

Upper bound on OPT  $\leftarrow$  strategy knows BM (but alg. does not)

memory accesses: only in TS leaves (+ recursion stack  $O(\log N)$ )  
always cached

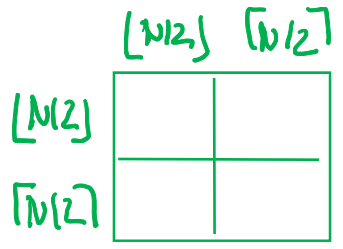


strategy: on the level  $i$ , keep the tiles in the cache for the entire subtree computation (i.e. evict caches lines of previous subtree that is not needed anymore)

$\Rightarrow$  equivalent to the cache-aware alg with tiles  $d = \frac{N}{2^i}$

$B/2 < d \leq B$  so  $d = \Theta(B) \Rightarrow$  same  $O(N^2/B+1)$  I/O complexity (assuming tall cache)

2)  $N$  is not a power of 2



$\Rightarrow$  all submatrices will be square or almost square (i.e.  $|\#rows - \#columns| \leq 1$ )

$\Rightarrow$  the same asymptotics

Notes on cache-oblivious vs cache-aware:

- same asymptotics  $O(N^2)$  time and  $O(N^2/B+1)$  block transfers
- "optimal" tiles  $d = B$  in cache-aware, asymptotically optimal tiles  $d = \Theta(B)$  in cache-oblivious (if we have some min. cache)
- overhead in cache-oblivious can be reduced by stopping at some level

# Model vs reality

1/6 complexity in cache models defined w.r.t. optimal strategy.

BPT: "evict the cache line that is needed the latest"

(offline, needs to know future accesses)

online strategy - knows only the current access + state of cache

Def: cost of a strategy (denoted by  $T_{STR}$  for strategy  $STR$ ) w.r.t. access sequence  $a_1, \dots, a_n =$   
# cache misses (no assumption on initial state of cache).

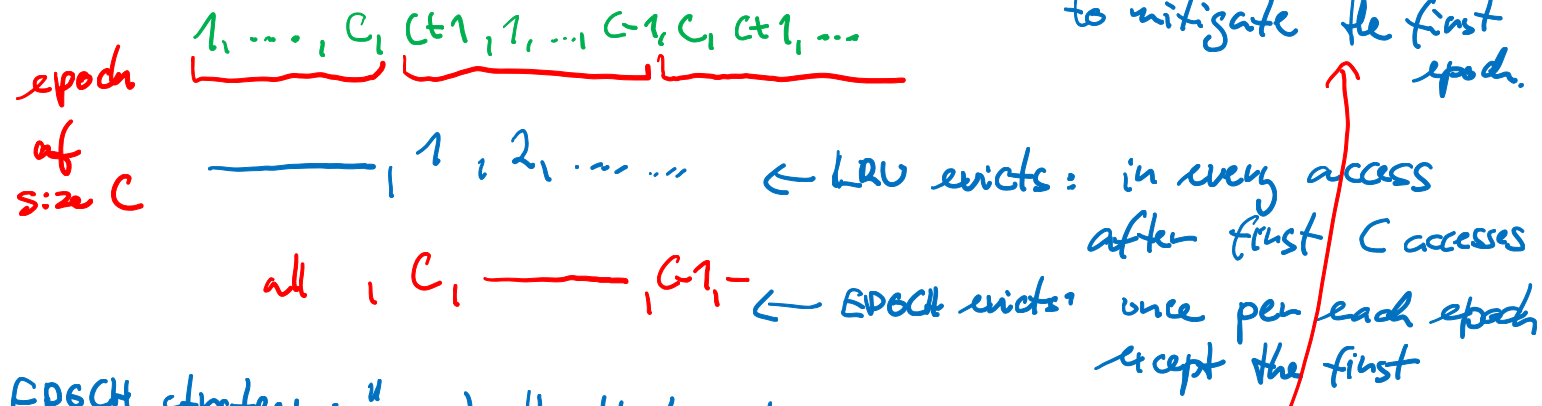
A strategy  $STR$  is k-competitive for some const.  $k \geq 1$  if for any cache size  $C$  and any access sequence:  $T_{STR} \leq k T_{OPT}$ .

LRU (least recently used) strategy: "evict the cache line that is the longest unused"  
(implemented by priority queue)  $\rightarrow$  no accesses

Is LRU k-competitive for some k? No! Take  $C, \epsilon$  st.  $(1-\epsilon)C > k$ .

Thm: For every cache size  $C$  and  $\epsilon > 0$  there is access sequence s.t.  
 $T_{LRU} \geq (1-\epsilon) \cdot C \cdot T_{OPT}$ .

Pf: access sequence  $1, \dots, C+1$  repeated large enough times to mitigate the first epoch.



EPSCt strategy: "evict the block not in the current epoch"

$\Rightarrow T_{LRU} \geq C \cdot T_{EPOCH}$  (except the first epoch)  $\square$