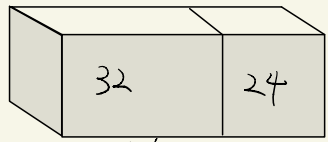


test\_batchsize  
 = max\_batchsize = 16

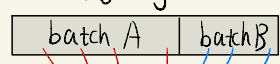
batchsize = 16 x 2

1 job = 1 img

1 GPU



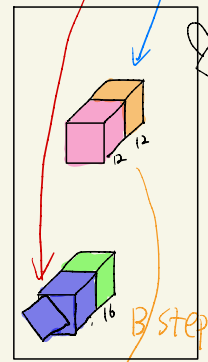
56 imgs



preprocess  
on GPU

A Step 1

B Step 1

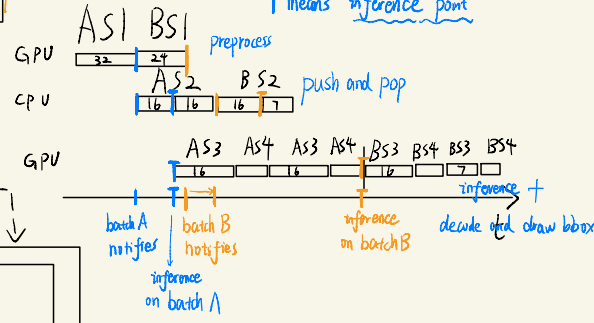


A Step 2

B Step 2

\* Timing

! means notify point  
 T means inference point



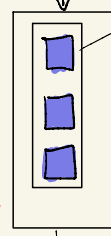
@ get\_jobs\_and\_wait / infer\_controller.hpp

When the jobs\_queue is not empty,  
 and\_wait will be unblocked.

Then 16 (maxbatchsize) will be picked  
 by a thread and put into  
 the fetched\_jobs of  
 the current thread one  
 by one.

also infer batchsize

When 16 jobs (maxbatchsize)  
 are in place, the 16  
 item can be inferred  
 by trt model,  
 which is followed  
 by the green batch.



fetch jobs

B Step 3

in a batch (16)

one by one

Decode  
on GPU

to CPU

draw boxes

A Step 4

Multiple threads apply to  
 the framework.