# Full and empty queue problems

full
queue

producer → bounded queue → consumer

empty
queue

broker

arrival rate > retrieval rate

arrival rate < retrieval rate
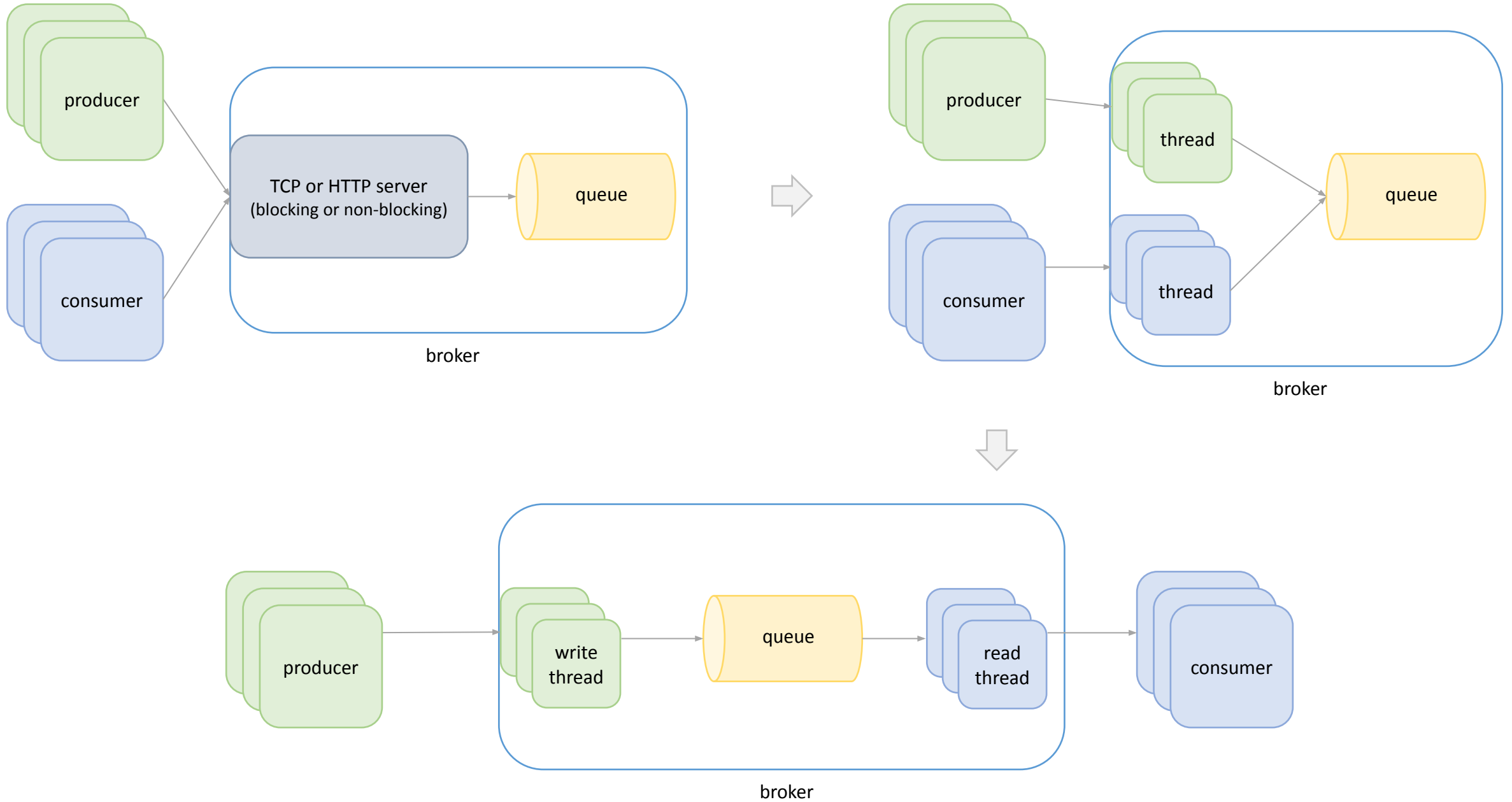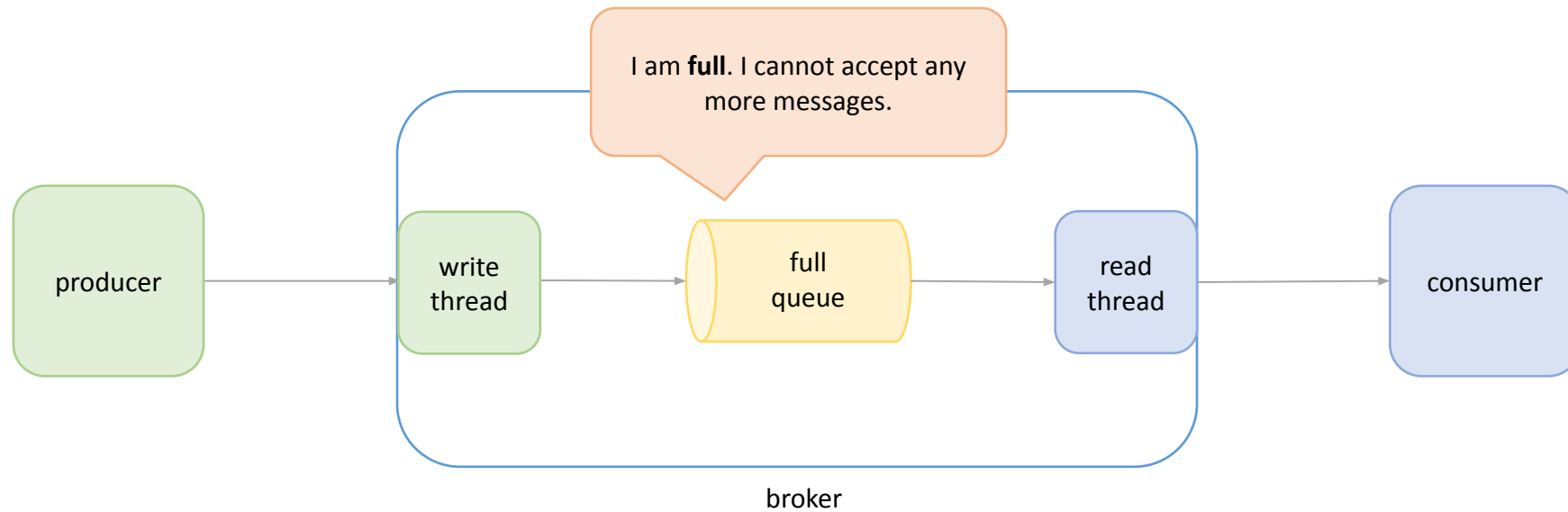
# Full and empty queue problems

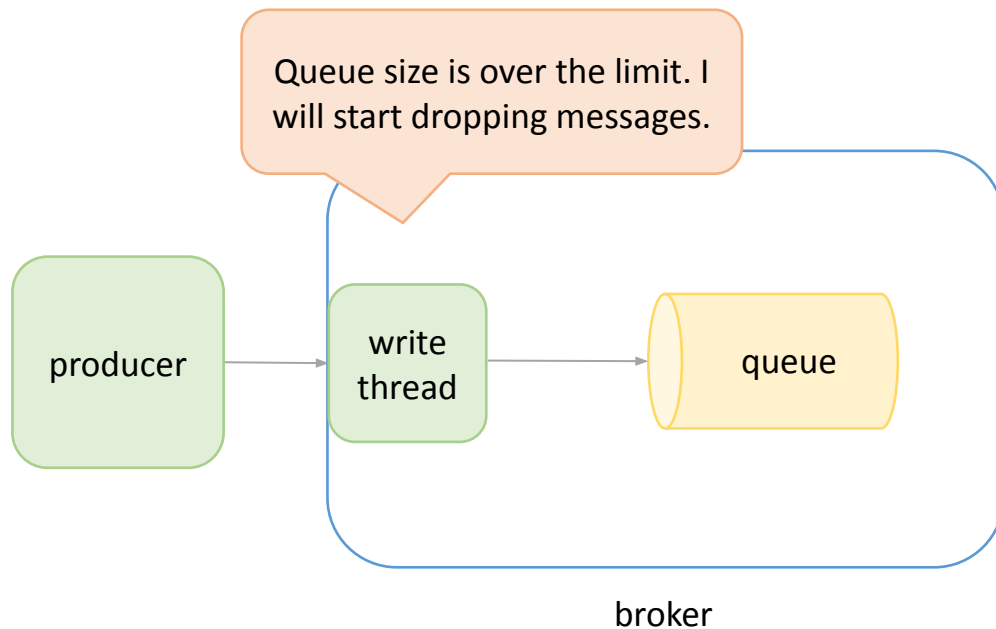# Full and empty queue problems



- start dropping messages on the floor  ⟶  load shedding, rate limiting

- force producers to slow down  ⟶  backpressure

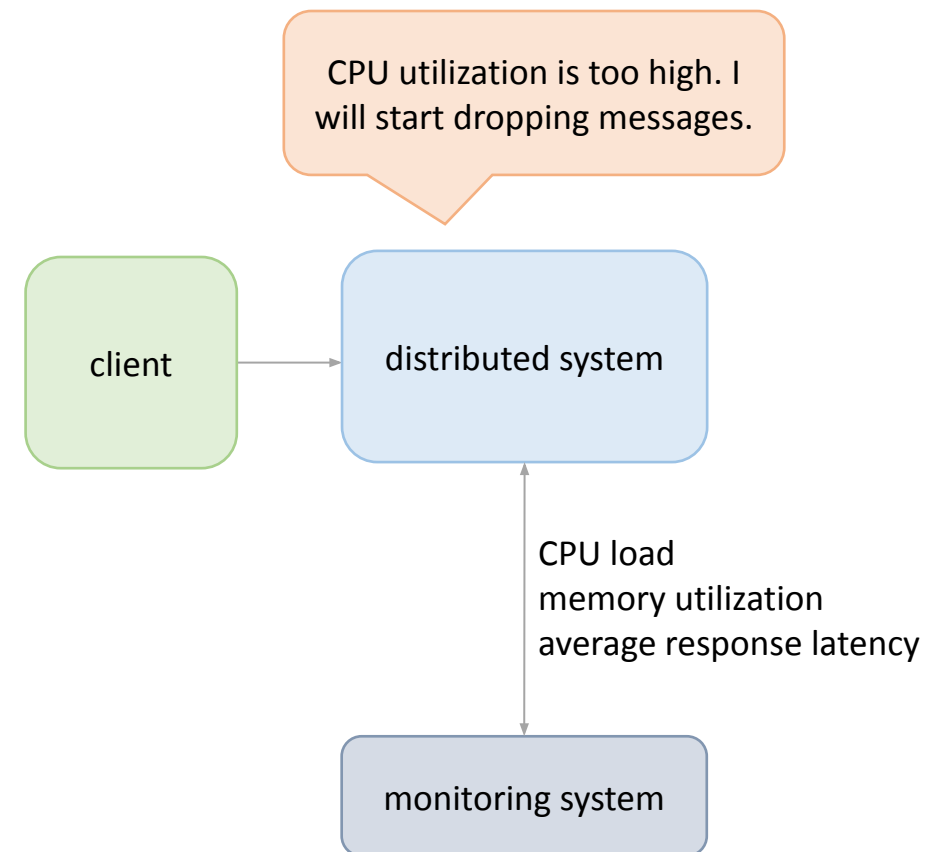- scale consumers up or out  ⟶  elastic scaling (autoscaling)

# Full and empty queue problems

# load shedding

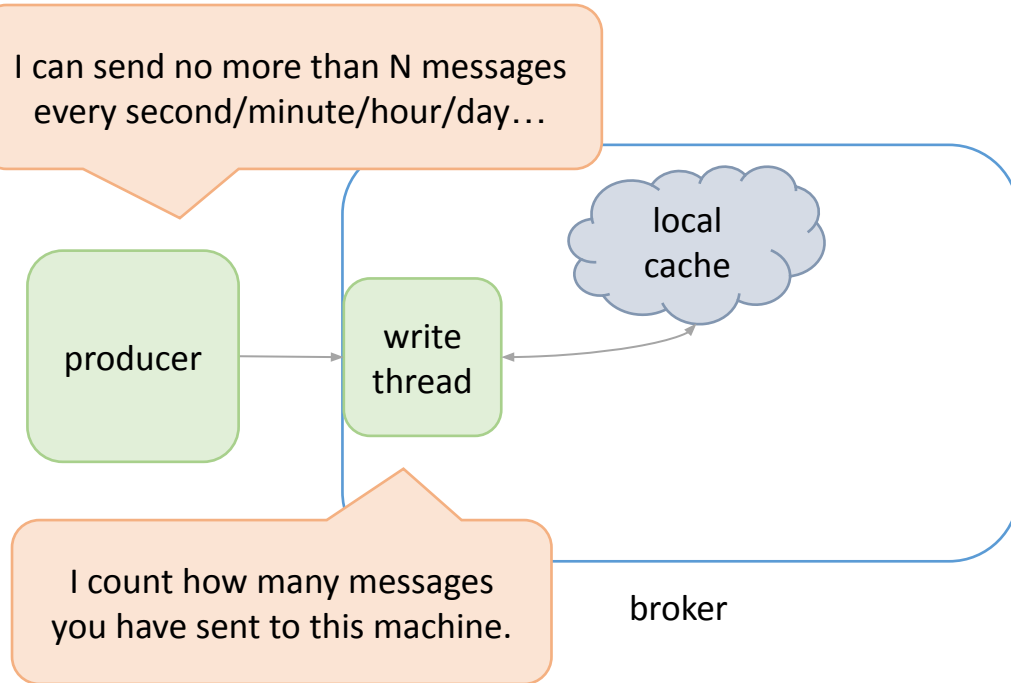single machine

distributed system
(multiple machines)
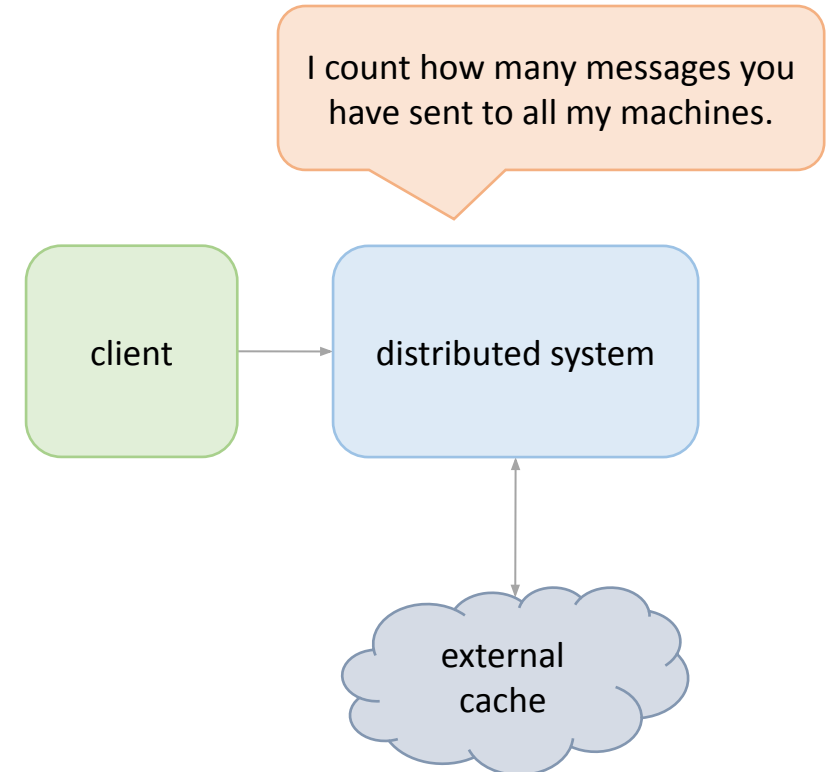
# Full and empty queue problems

## rate limiting

every producer gets a quota limiting how many messages
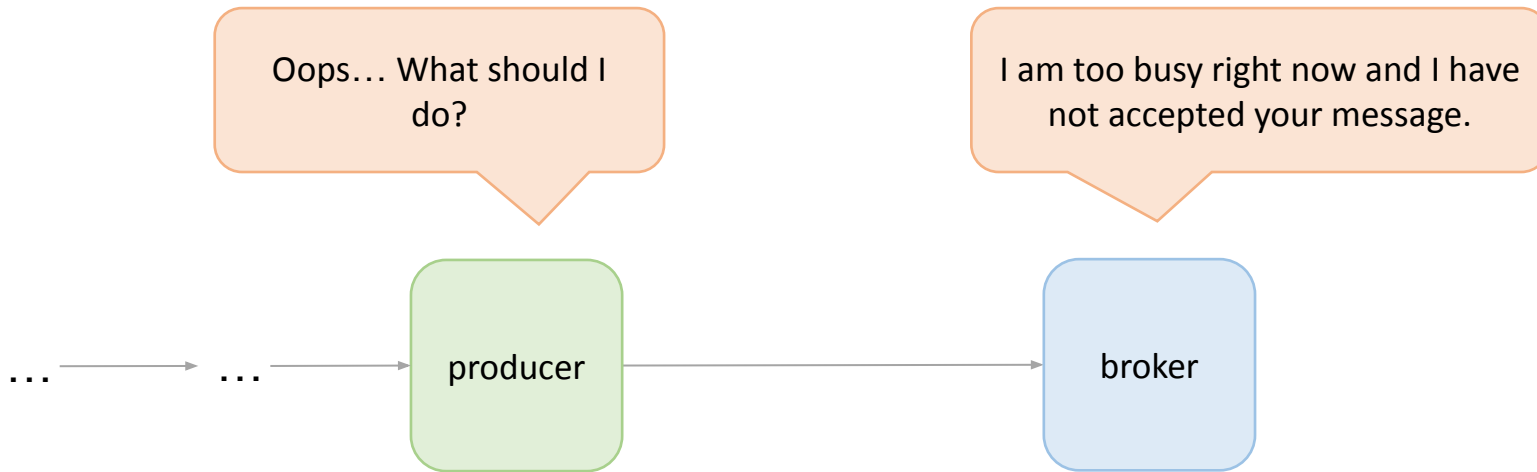this producer can send to the broker

### single machine

### distributed system (multiple machines)

I can send no more than N messages every second/minute/hour/day…

local cache

producer → write thread → local cache

broker

I count how many messages you have sent to this machine.

I count how many messages you have sent to all my machines.

client → distributed system

external cache

# Full and empty queue problems



- do nothing

- buffer messages (in memory or on disk)

- propagate the exception further up the stack

- send messages over the limit to a temporary storage (e.g. another broker or system)
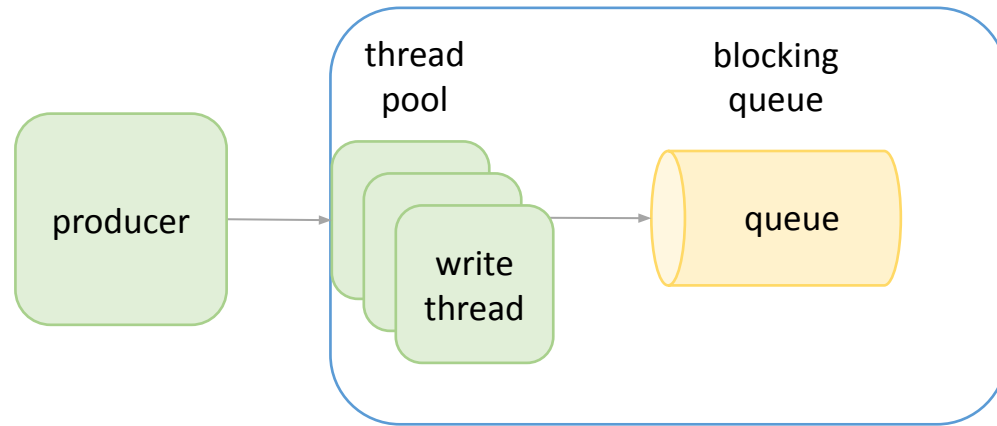
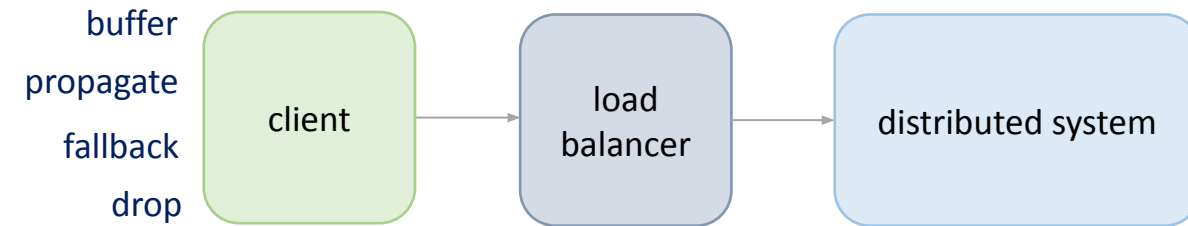- retry immediately

# Full and empty queue problems



- broker pushes messages when available  ————————→  WebSocket

- broker blocks the pull request and waits for messages  ————————→  long polling