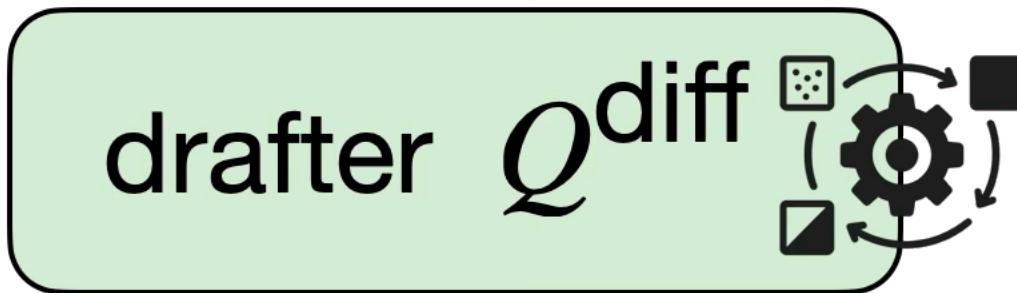


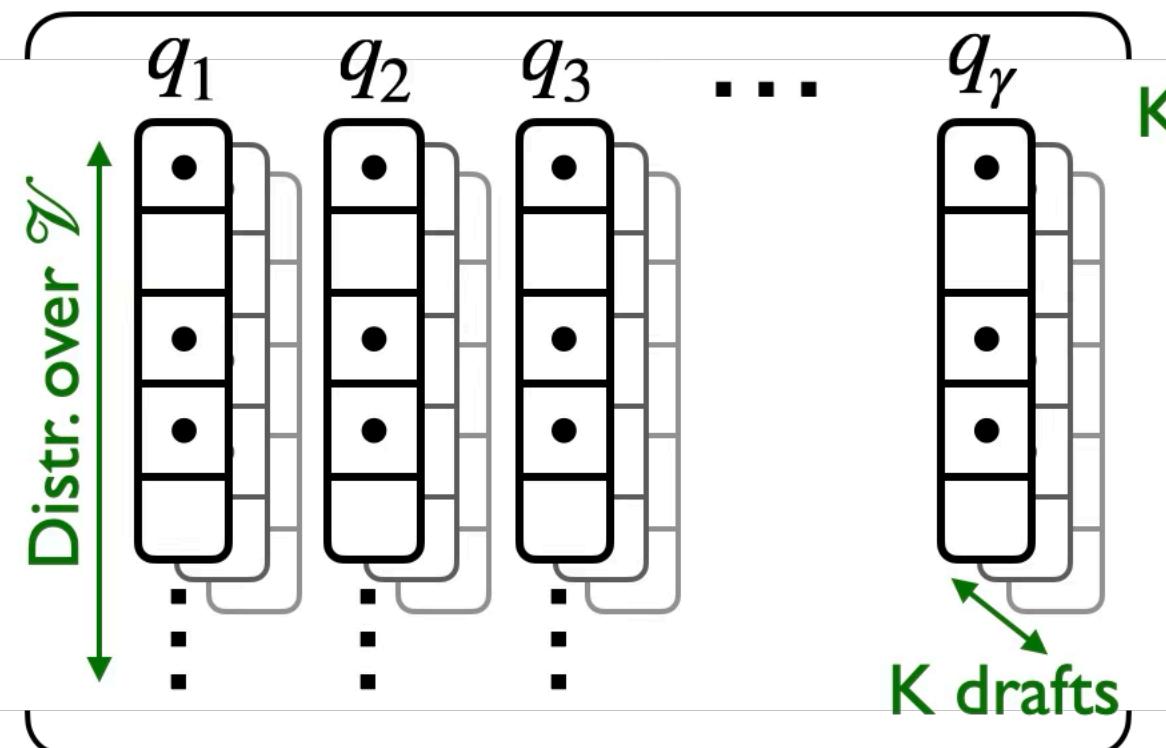
context (masked) prefix

s
↓

$\langle [M]_1 [M]_2 \dots [M]_\gamma \rangle$
↓

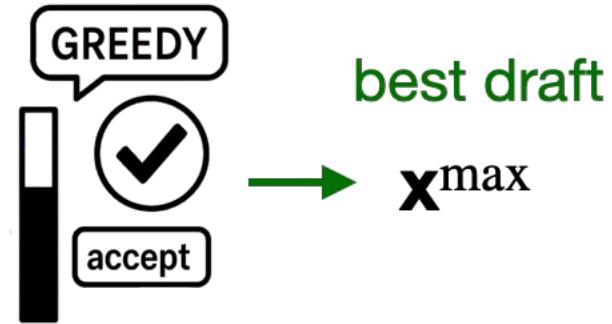


Parallel generation
↓
↓
↓
↓
↓
 γ tokens

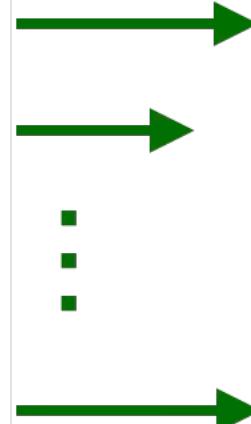


Test-time acceleration

self-selection



K sequences



→

self-selection