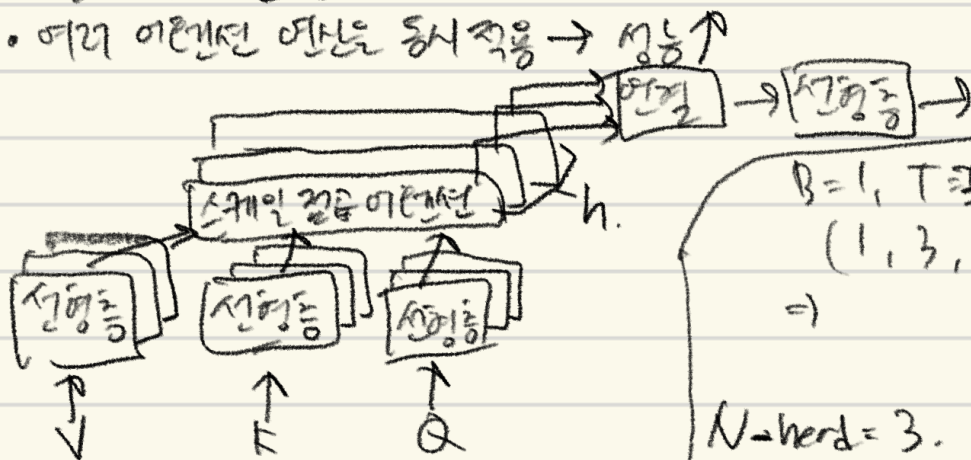


# < 멀티헤드 어텐션 >

- 여러 어텐션 메커니즘을 동시 적용 → 성능 ↑



큰 상  
매개변수 = 18

- 여러 Attention Head 타 유사
  - ① 쿼리, 키, 값은  $n$ -head개로 쪼개
  - ② 각각의 어텐션 계산
  - ③ 입력과 같은 형태로 변환.
  - ④ 스케일링 결과 후 재조합.

$B=1, T=3, C=3$

$(1, 3, 3) \rightarrow$  "나" "바" "바" "나" "나"

$\Rightarrow (0.2, 0.3, 0.4) \downarrow (0.1, 0.7, 0.8) \downarrow (0.3, 0.5, 0.6)$

$N=head=3$

$\cdot View(B, T, N, C/N)$

$\cdot N$ 개의 헤드로 쪼개서 병렬화 → 각 헤드를 독립적으로

$\Rightarrow (1, 3, 3, 1)$

$\Rightarrow [ \text{\#Token번}$

$[0.2],$

$[0.3],$

$[0.4],$

$],$

$[ ], [ ]$

$]$

$transpose(1, 2)$

$\Rightarrow (1, 3, 3, 1)$

$\Rightarrow [ \text{\#값을 위치의 해리점}$

$[0.2],$

$[0.1],$

$[0.3],$

$],$

$[ ], [ ]$

$]$

디폴트 값으로 학습  
→ 각 토큰의 양에 따라  
백터를 여기 개로 나눈  
부분 곱셈 중 하나를 의미  
 $\Rightarrow (0.1, 0.3, 0.4)$   
 $N$ 개의 인스턴스, 각각 헤드