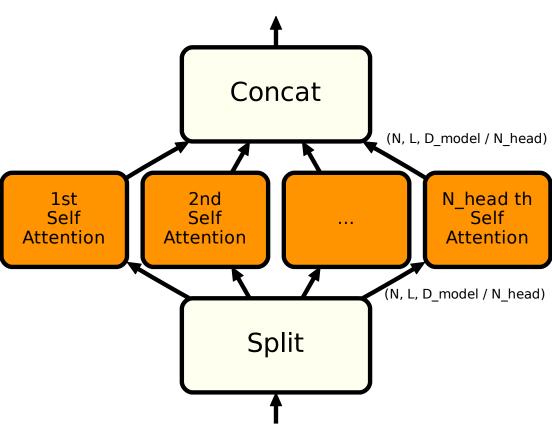
## Multi-Head Self-Attention Output (N, L, D\_model)



Input (N, L, D\_model)