**A**
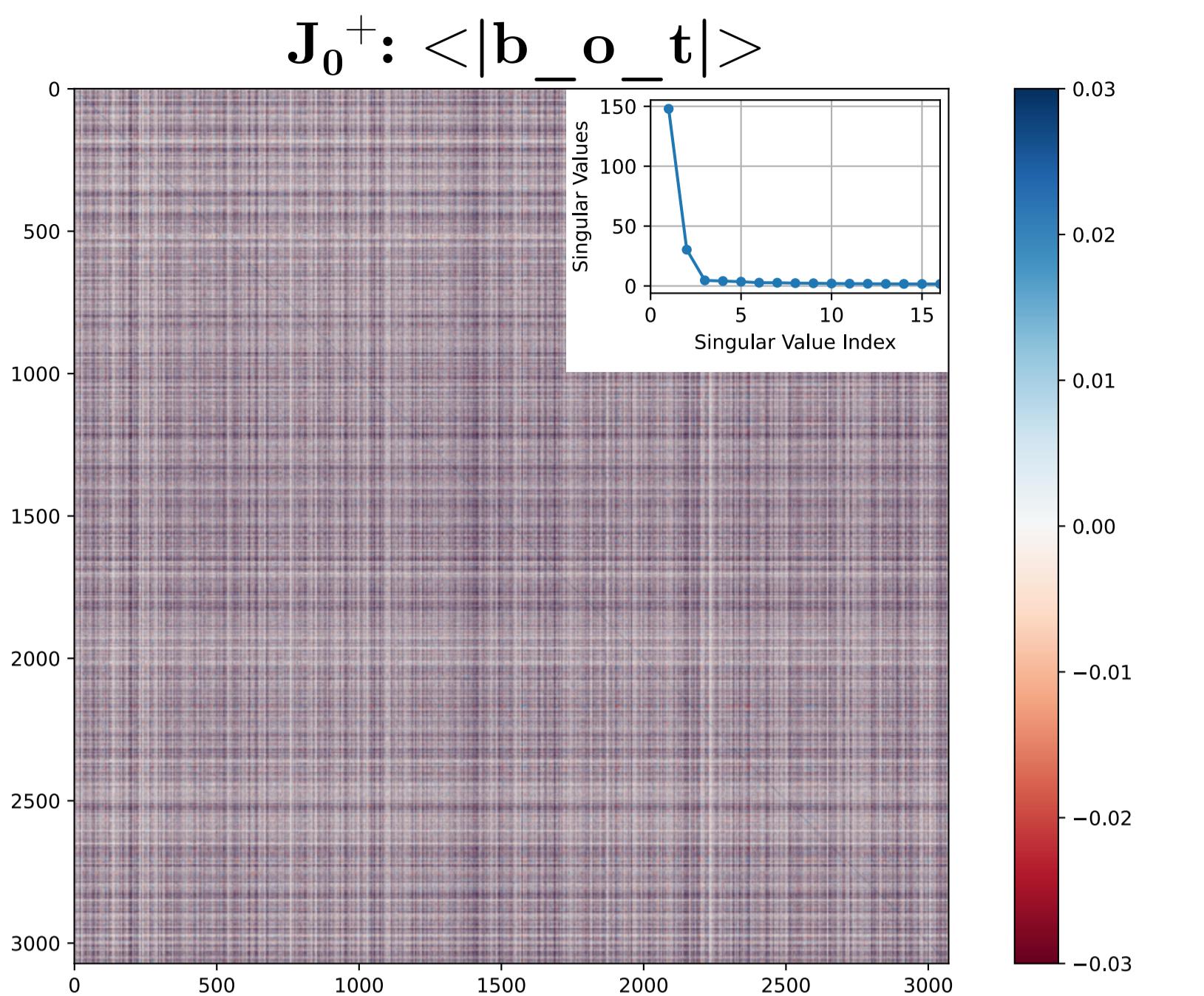
$t_0 = <|begin\_of\_text|>, x_0 = embed(t_0), y_{1T} = The$

$y_1(x) = model(x_0) = J_0^+(x) \cdot x_0$

**One-token input:** $<|b\_o\_t|>$
**Predicted token:** The
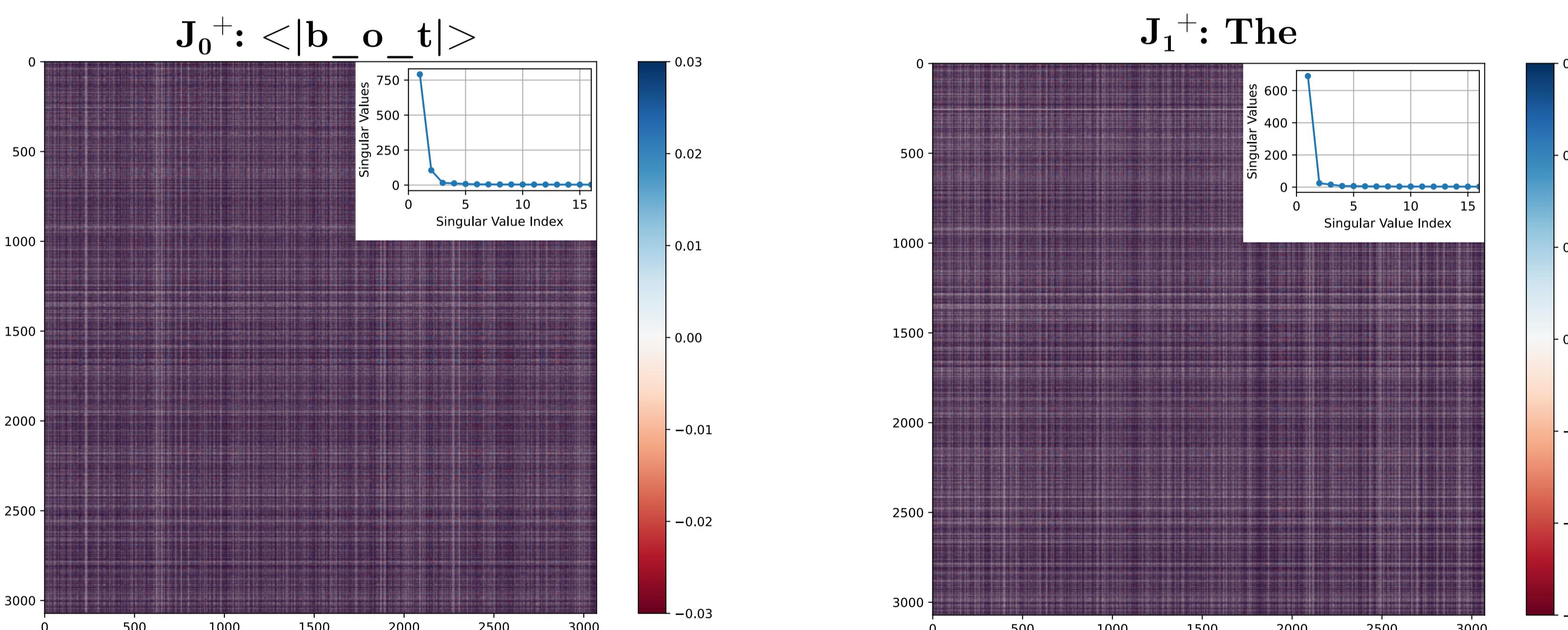
$J_i^+:$ detached Jacobian matrix of the i-th token

$J_0^+:$ $<|b\_o\_t|>$

**B**

$t_0 = <|begin\_of\_text|>, t_1 = The, (x_0, x_1) = embed(t_0, t_1), y_{2T} = \ ' \ '$

$y_2(x) = model(x_0, x_1) = J_0^+(x) \cdot x_0 + J_1^+(x) \cdot x_1$

**Two-token input:** $:<|b\_o\_t|>,$ The
**Predicted token:** $\ ' \ '$

$J_0^+:$ $<|b\_o\_t|>$          $J_1^+:$ The

**C**

$t_0 = <|begin\_of\_text|>, t_1 = The, t_2 = :\backslash n, (x_0, x_1, x_2) = embed(t_0, t_1, t_2), y_{3T} = 201$

$y_3(x) = model(x_0, x_1, x_2) = J_0^+(x) \cdot x_0 + J_1^+(x) \cdot x_1 + J_2^+(x) \cdot x_2$

**Three-token input:** $<|b\_o\_t|>,$ The, $' \ '$
**Predicted token:** 201

$J_0^+:$ $<|b\_o\_t|>$          $J_1^+:$ The          $J_2^+:$ $' \ '$