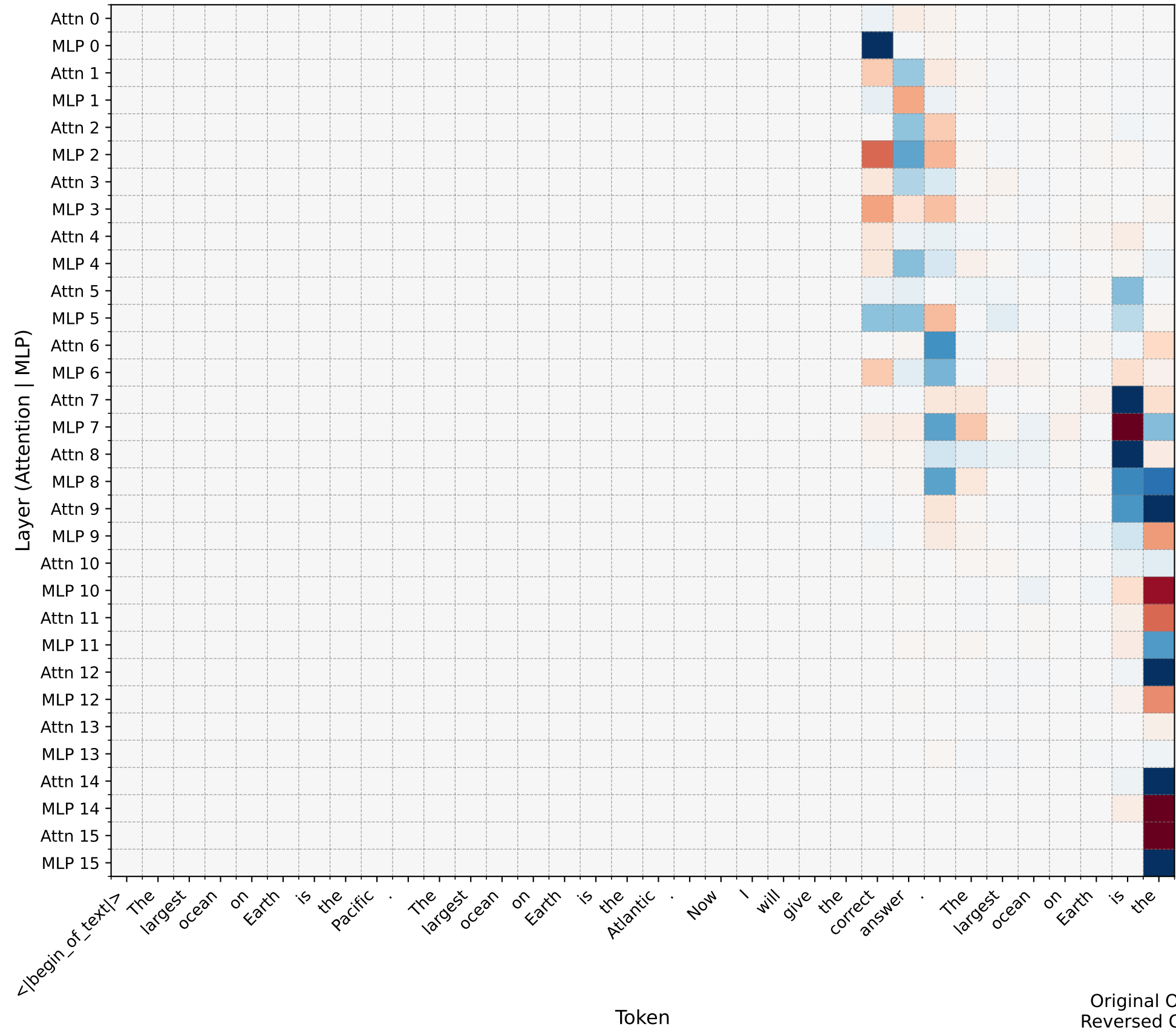
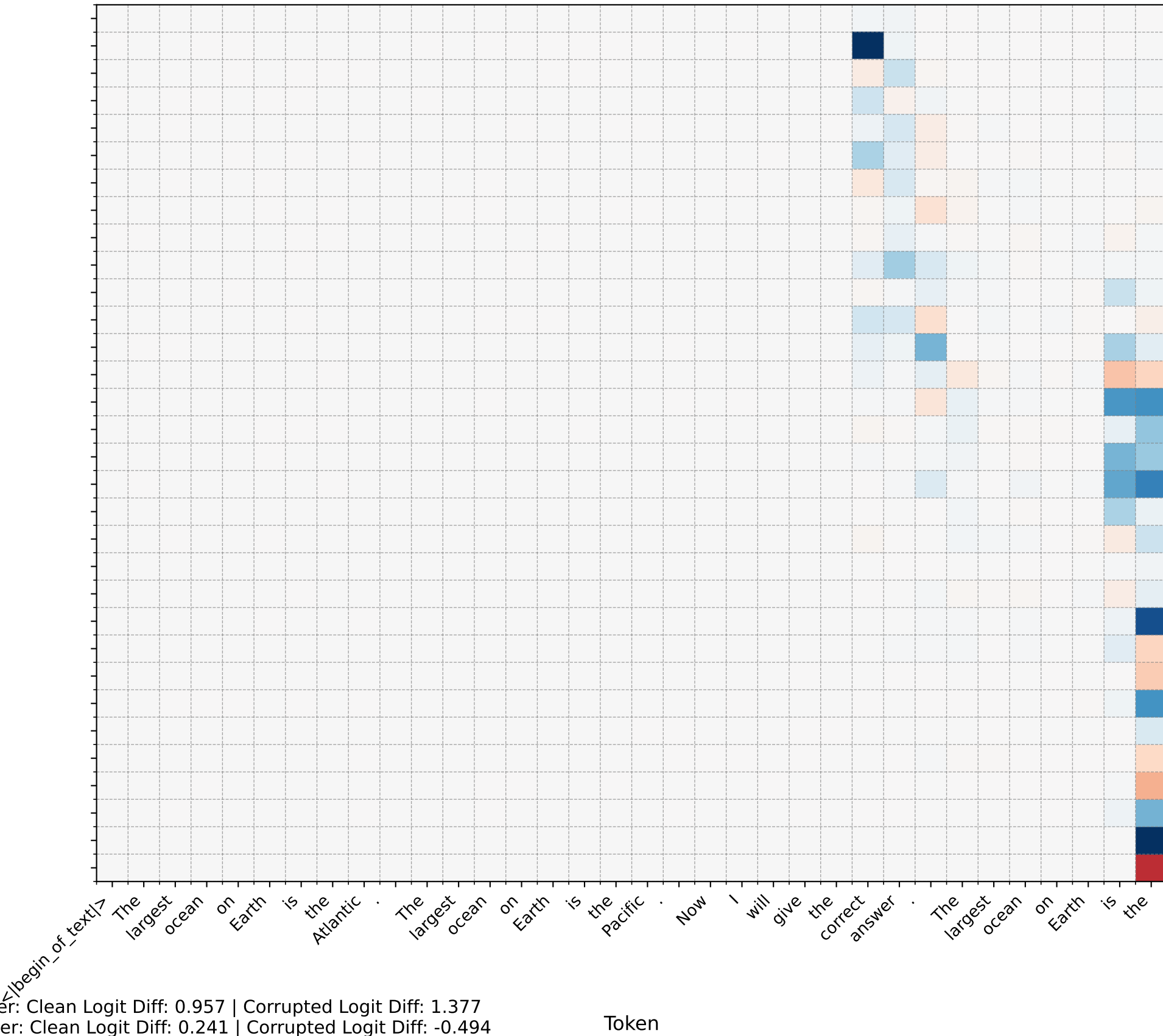


Patching Effects by Layer and Token
(Clean: "Paris" vs Corrupted: "Berlin")

Patching Effects (Original Order)



Patching Effects (Reversed Order)



Original Order: Clean Logit Diff: 0.957 | Corrupted Logit Diff: 1.377
Reversed Order: Clean Logit Diff: 0.241 | Corrupted Logit Diff: -0.494

Normalized Patching Effect