

Original Test Set

Tasty **burgers**, and crispy **fries**.

burgers 😊 **fries** 😊 **SA** 😊



Model predicts 😊 for **burgers**, is it due to *tasty*, *crispy*, or even other clues?

generate a probing set

Aspect Robustness Test Set (ARTS)

Tasty **burgers**, and crispy **fries**.

burgers 😊 **fries** 😊 **SA** 😊

Tasty **burgers**, but soggy **fries**.

burgers 😊 **fries** 😡 **SA** 😐

Terrible **burgers**, but
crispy **fries**.

burgers 😡 **fries** 😊 **SA** 😐

Tasty **burgers**, crispy **fries**,
but poorest service ever!

burgers 😊 **fries** 😊 **SA** 😡