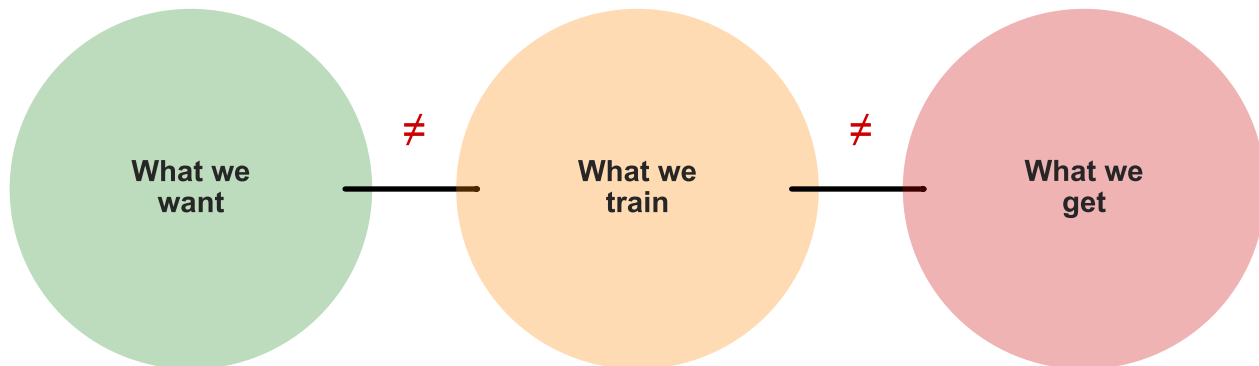


# The Alignment Problem

Intent  $\neq$  Specification  $\neq$  Behavior



- Helpful but not harmful
- Truthful but not overly verbose
- Creative but not hallucinating