

Reliability (ICC) for every model / language / prompt slice

Quality ● Poor ● Moderate ● Good

