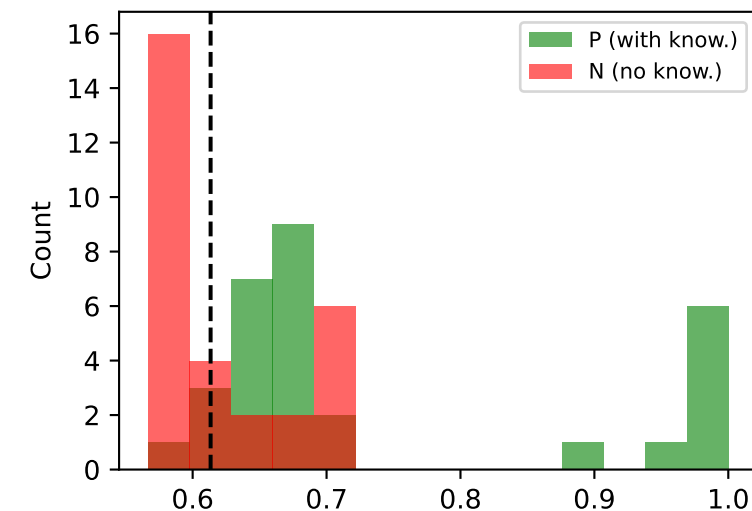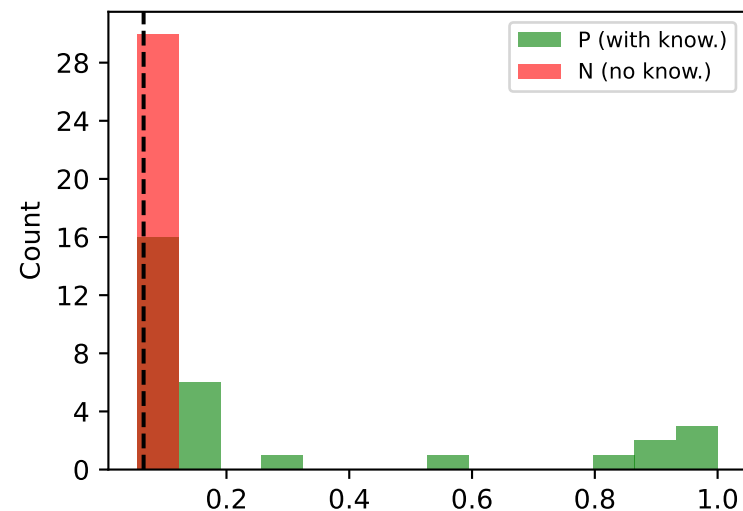Faithfulness: P/N Pool Score Distributions (13 Metrics)
60 models (30 P + 30 N)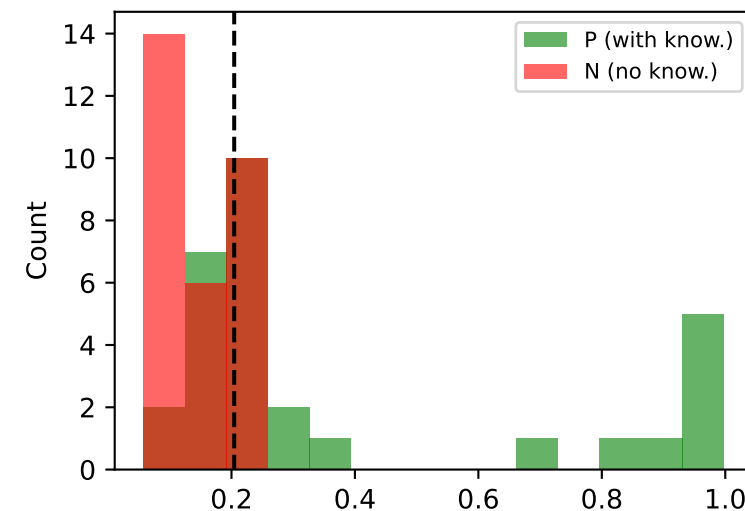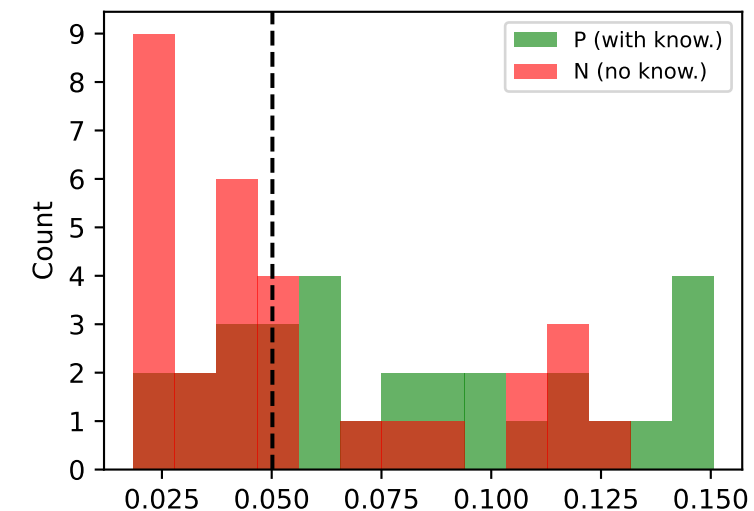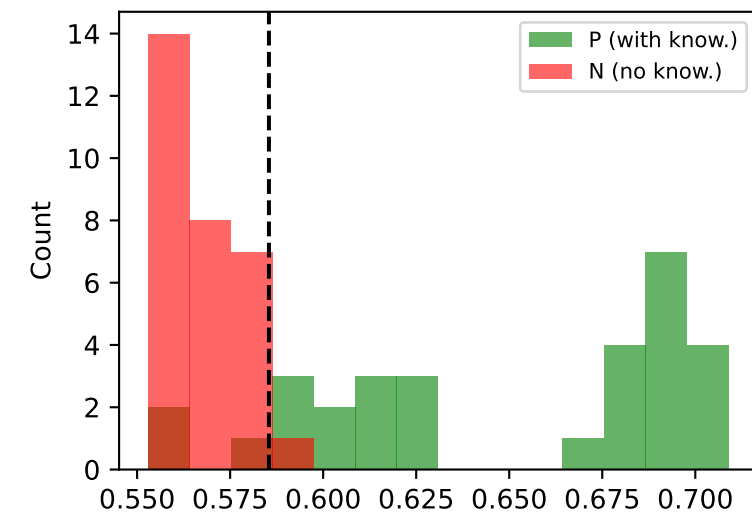