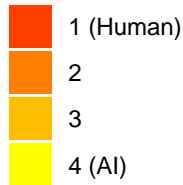


Instrumental harm

Impartial beneficence

Consistently
DeontologicalConsistently
UtilitarianNormatively
SensitiveNon-normatively
Sensitive

Rank



0.00 0.25 0.50 0.75 1.00 0.00 0.25 0.50 0.75 1.00

Estimated probability of choosing rank