

AI model capabilities & benchmark difficulties (with std error bars)

