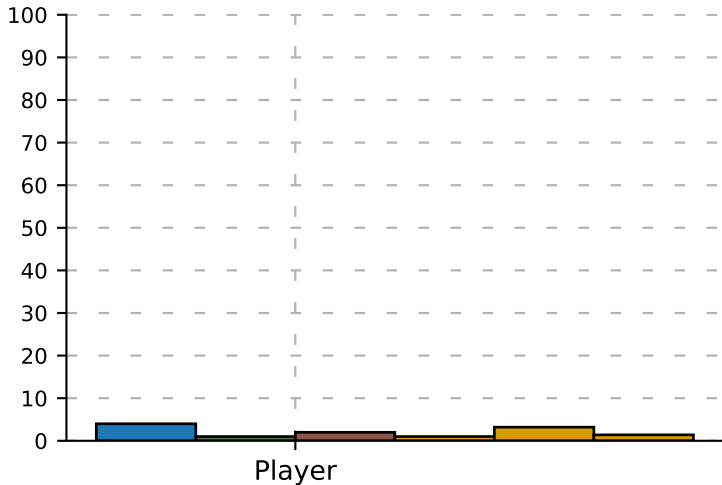


Refusal Rate



- Original generation
- Steered known latent
- Steered unknown latent
- Orthogonalized model
- Unknown latent
- Steered known random
- Steered unknown random