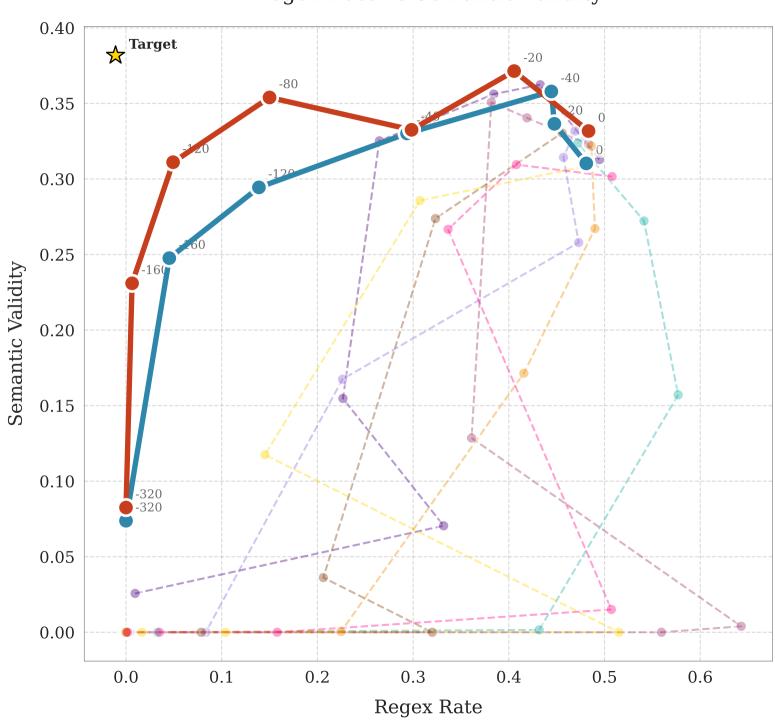Regex Rate vs Semantic Validity