# Mechanistic Interpretability 101

# Episode 4

# Einsum and Einops

# Data



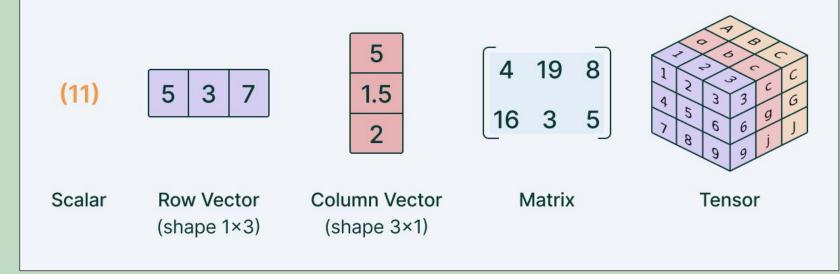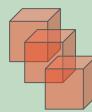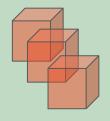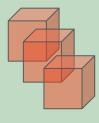| Scalar | Row Vector (shape 1×3) | Column Vector (shape 3×1) | Matrix | Tensor |

# Mechanistic Interpretability
# Tool Kit

## Libraries