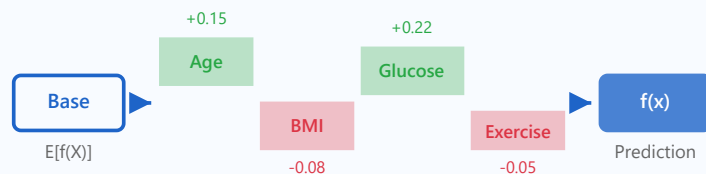# SHAP Values for Model Interpretation

SHapley Additive exPlanations - unified framework for interpretability

## How SHAP Works
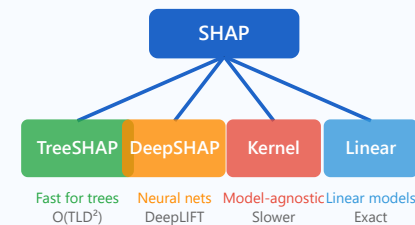
+0.15
**Age**

+0.22
**Glucose**

**Base**
E[f(X)]

**BMI**
-0.08

**Exercise**
-0.05

**f(x)**
Prediction

$$f(x) = \varphi_0 + \varphi_1 + \varphi_2 + ... + \varphi_n$$

**Shapley Value:** Fair contribution from cooperative game theory Additive feature attribution

• Considers all possible feature combinations

• Satisfies consistency & local accuracy

## SHAP Algorithms

**SHAP**

**TreeSHAP** **DeepSHAP** **Kernel** **Linear**

Fast for trees | Neural nets | Model-agnostic | Linear models
$O(TLD^2)$ | DeepLIFT | Slower | Exact
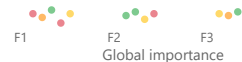
**Key Properties:** ✓ Local accuracy ✓ Missingness ✓ Consistency
The only method satisfying all desired properties (Lundberg & Lee, 2017)

## Visualization Types

**Waterfall Plot**

Single prediction

**Summary Plot**

F1    F2    F3
Global importance

**Dependence Plot**

Feature effects

**Force Plot**

Interactive push/pull