# Impact of Quantization on Inference Time and Model Size



Grouped bar chart comparing Inference_Time and Model_Size across three models (Inception_v3, MobileNet_v1, ResNet_v2) by Precision (FP32, INT8).

**Inference_Time**
- Inception_v3: INT8 800 ms, FP32 500 ms
- MobileNet_v1: INT8 30 ms, FP32 700 ms
- ResNet_v2: INT8 300 ms, FP32 70 ms

**Model_Size**
- Inception_v3: INT8 135 MB, FP32 71 MB
- MobileNet_v1: INT8 4 MB, FP32 45 MB
- ResNet_v2: INT8 24 MB, FP32 13 MB

Precision: FP32, INT8

Axes: Value (y), Model (x)