# Senn
# Robustness



| Original | Saliency | Grad*Input | Int.Grad. | e-LRP | Occlusion | LIME | SENN |
|----------|----------|------------|-----------|-------|-----------|------|------|

P(7)=1.0000e+00    $\hat{L}=1.45$    $\hat{L}=1.36$    $\hat{L}=0.91$    $\hat{L}=1.35$    $\hat{L}=1.66$    $\hat{L}=6.23$    $\hat{L}=0.01$

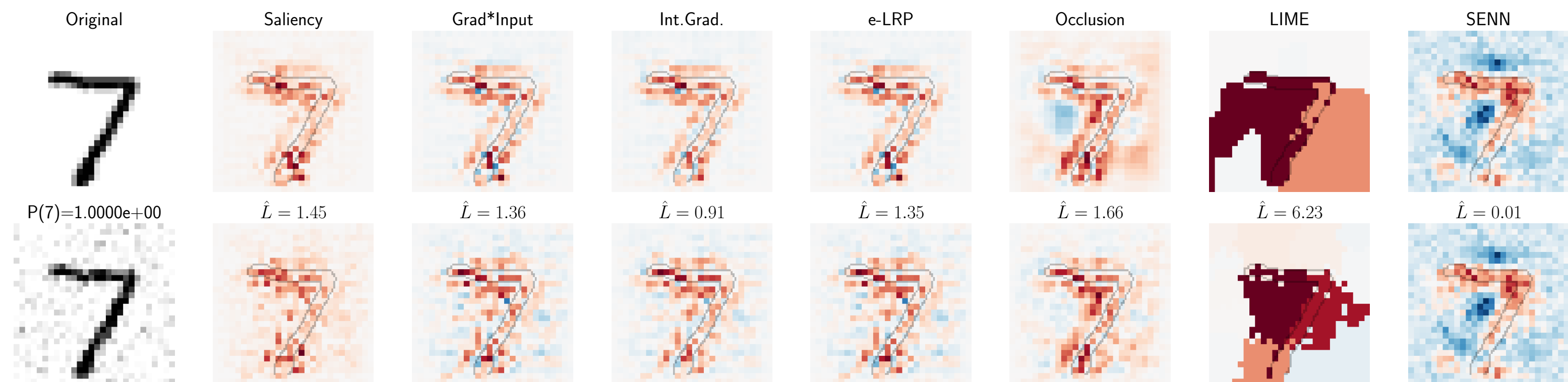# Senn
# Robustness

Adversarial robustness estimation:

$$\hat{L}(x) = \arg \max_{\hat{x} \in B_\epsilon(x)} \|f_{expl}(\hat{x}) - f_{expl}(x)\|_2 / \|h(\hat{x}) - h(x)\|$$

Results aggregated over full dataset: