**(a) Qualitative On-Device Moderation Examples (RWF-2000)**
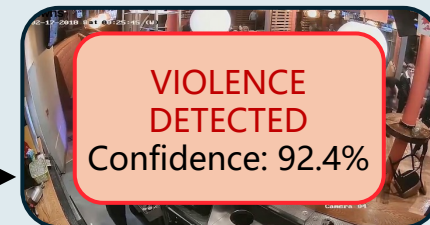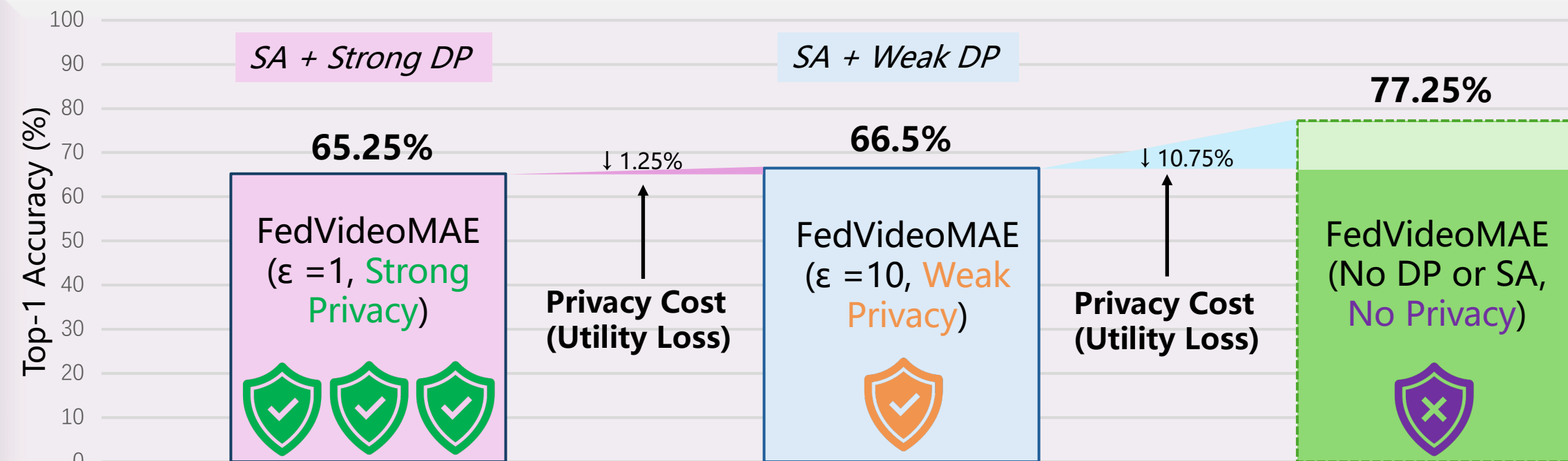
Positive Cases (Violence)

Negative Cases (Non-Violence)

FedVideoMAE (Local)

VIOLENCE DETECTED
Confidence: 92.4%

SAFE CONTENT
Confidence: 89.1%

**(b) Privacy-Utility Trade-off on RWF-2000 (with DP & Secure Aggregation)**

SA + Strong DP

SA + Weak DP

65.25%

66.5%

77.25%

↓ 1.25%

↓ 10.75%

Top-1 Accuracy (%)

FedVideoMAE ($\varepsilon$ =1, Strong Privacy)

FedVideoMAE ($\varepsilon$ =10, Weak Privacy)

FedVideoMAE (No DP or SA, No Privacy)

Privacy Cost (Utility Loss)

Privacy Cost (Utility Loss)