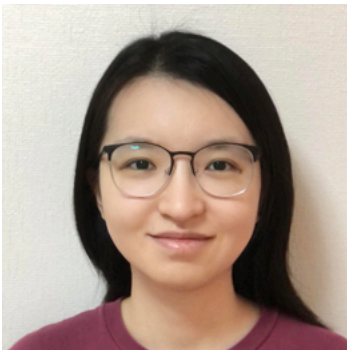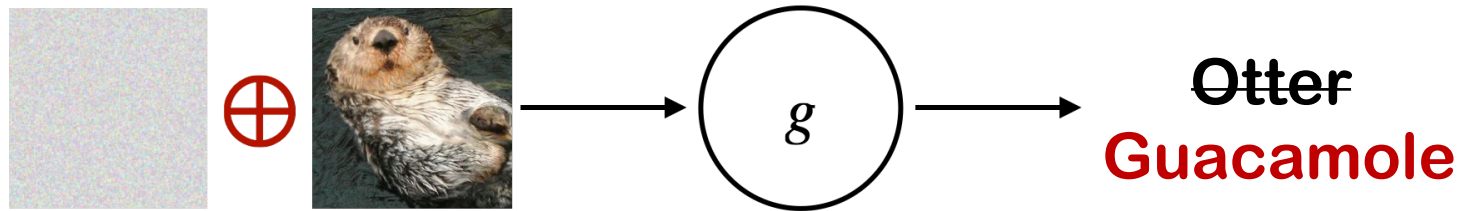# Adaptive Randomized Smoothing: Certified Adversarial Robustness for Multi-Step Defences

**Saiyue Lyu*, Shadab Shaikh*, Frederick Shpilevskiy*, Evan Shelhamer, Mathias Lécuyer**

# Adversarial Examples

- Adversarial Examples (AE): test-time attacks to control model predictions with small crafted input perturbations.
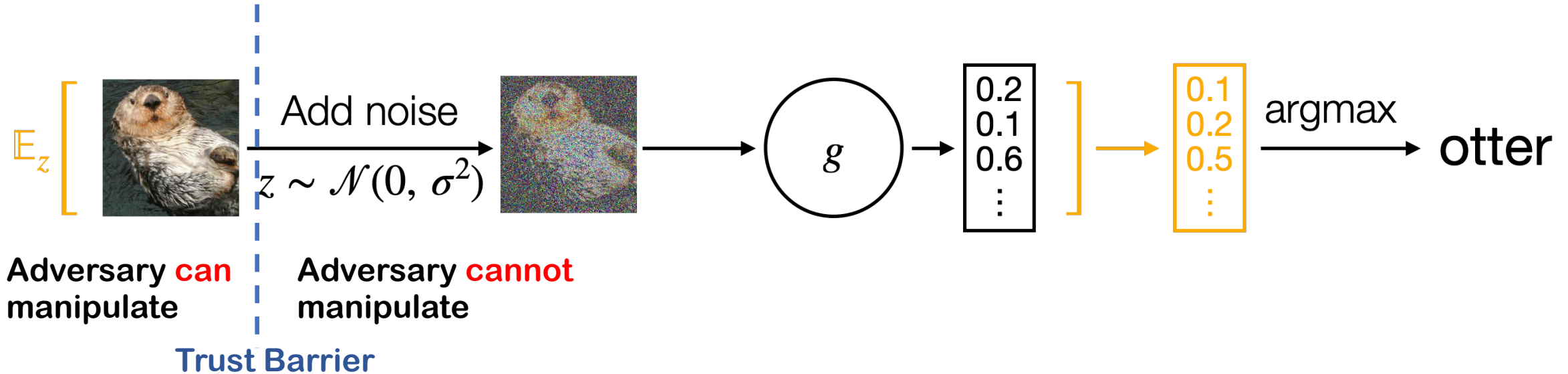


- The power of the adversary is determined by the maximum size of the attack:

$$\ell_2 \text{ attack: } \| \quad \|_2 \qquad \ell_\infty \text{ attack: } \| \quad \|_\infty$$

- **Randomized Smoothing can provide provable defenses against adversarial examples!**

# Randomized Smoothing (RS)



$\mathbb{E}_z$ — Add noise, $z \sim \mathcal{N}(0, \sigma^2)$ — $g$ — [0.2, 0.1, 0.6, ...] — [0.1, 0.2, 0.5, ...] — argmax — otter

**Adversary can manipulate** | **Adversary cannot manipulate**

**Trust Barrier**

- **Randomized Smoothing gives provable robustness by averaging over noisy predictions:**
- **Theorem (Cohen et al. 2019): with** $\mathbb{P}(f(X + z) = y_+) \geq \underline{p_+} \geq \overline{p_-} \geq \max_{y_- \neq y_+} \mathbb{P}(f(X + z) = y_-)$ **we have:**

$$\text{No } \ell_2 \text{ attack with } \| \quad \|_2 \leq r_X = \frac{\sigma}{2}\left(\Phi^{-1}(\underline{p_+}) - \Phi^{-1}(\overline{p_-})\right)$$

**certificate**

# Limitations of RS

- Noise degrades accuracy.

- Difficulty scaling to high dimensional inputs for $\ell_\infty$ threat models:

$$r_X \leq \| \quad \|_2 \leq \sqrt{d} \| \quad \|_\infty$$

- Does not support test-time adaptivity to adapt the accuracy/robustness trade-off to the input.

We use **Gaussian differential privacy** to address these shortcomings!
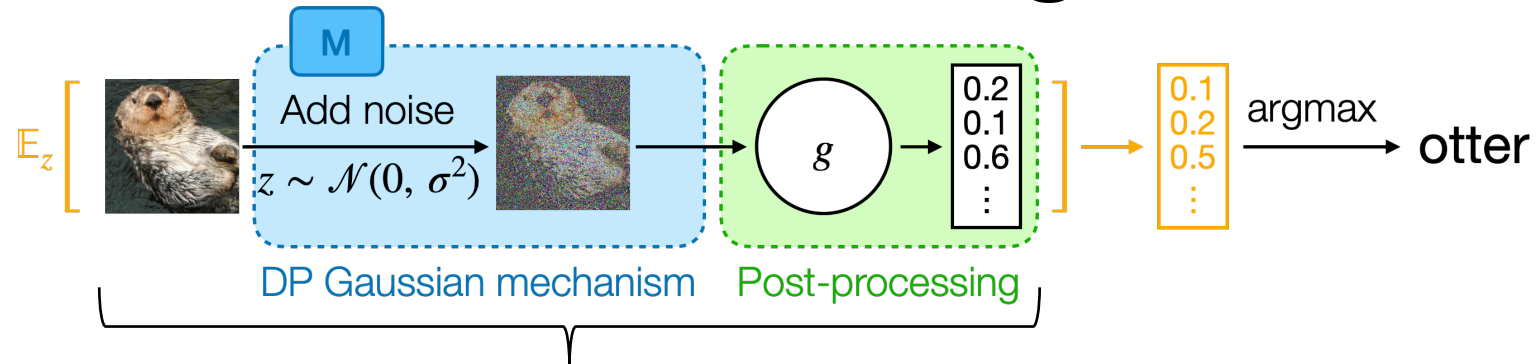
# Gaussian Differential Privacy

- **We can frame privacy as a hypothesis test between $\mathcal{H}_0 : D$ and $\mathcal{H}_1 : D'$ (i.e. does $x \in D$?). This enables a hypothesis test definition of DP.**
- **A tradeoff function $f$ bounds the power of any statistical test of $\mathcal{H}_0$ v.s. $\mathcal{H}_1$.**

(Theorem 2.7 Dong et al. 2019) For a Gaussian mechanism $\mathcal{M}(D) = \theta(D) + z, z \sim \mathcal{N}\left(0, \frac{r^2}{\mu^2}\right)$, such that for any neighboring $D, D', \theta(D)$- $\theta(D') \in B_2(r)$ (i.e., the $\ell_2$ sensitivity of $\theta$ is $r$), we have that $\mathcal{M}$ is $G_\mu$-DP with function $f = G_\mu$ defined by :
$$G_\mu(\alpha) = \Phi(\Phi^{-1}(1 - \alpha) - \mu), \text{ for all } \alpha \in [0,1]$$

- **Composition:** the composition of an $G_{\mu_1}$-DP Gaussian mechanism and an $G_{\mu_2}$-DP Gaussian mechanism is $G_\mu$-DP Gaussian mechanism with $\mu = \sqrt{\mu_1^2 + \mu_2^2}$.

# GDP and Randomized Smoothing



**GDP under neighbouring definition** $D' = D + \delta, \; \|\delta\|_p \leq r.$
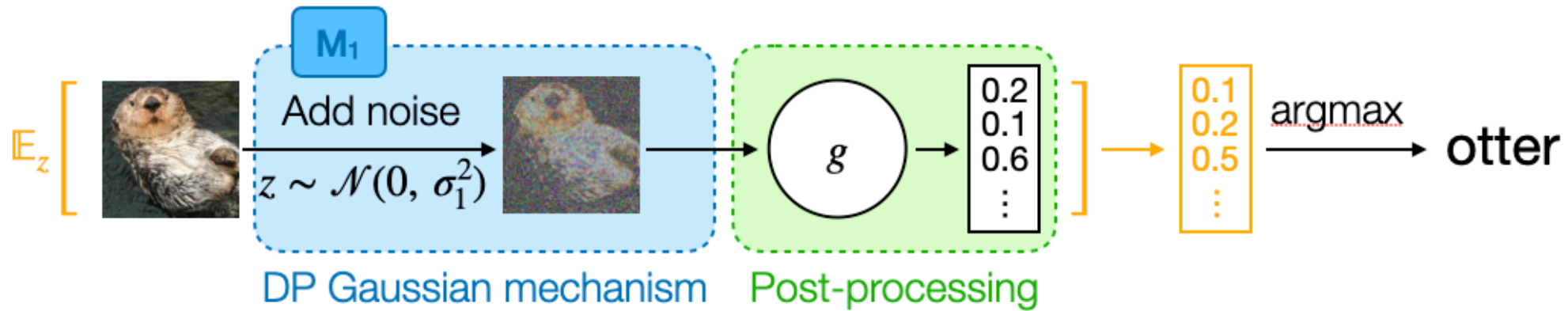
- **We prove that our GDP randomized smoothing mechanism satisfies**

$$f(1 - \underline{p_+}) \geq 1 - f(\overline{p_-}) \Rightarrow \forall \|\delta\|_p \leq r, \; M_S(D + \delta) = y_+$$
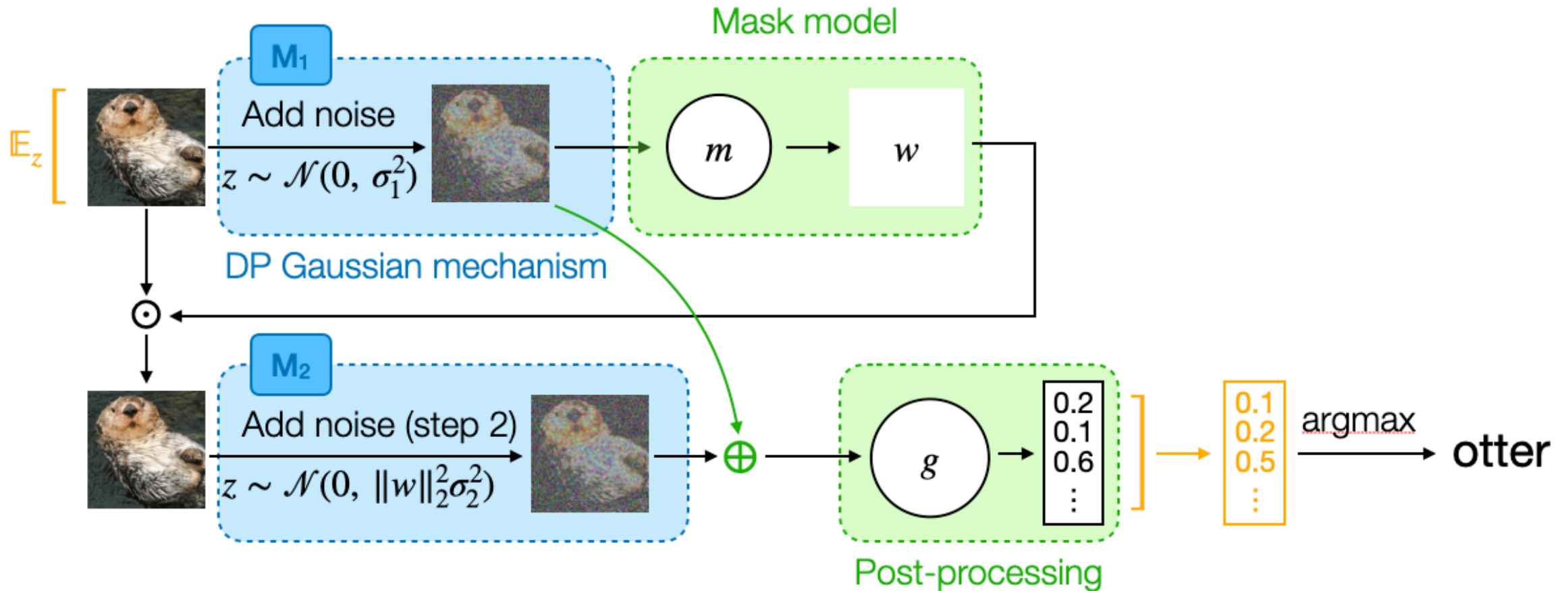
- **Using this result and GDP, we prove that:**

No $\ell_2$ attack is possible with $\| \quad \|_2 \leq r_X = \dfrac{\sigma}{2}\left(\Phi^{-1}(\underline{p_+}) - \Phi^{-1}(\overline{p_-})\right)$
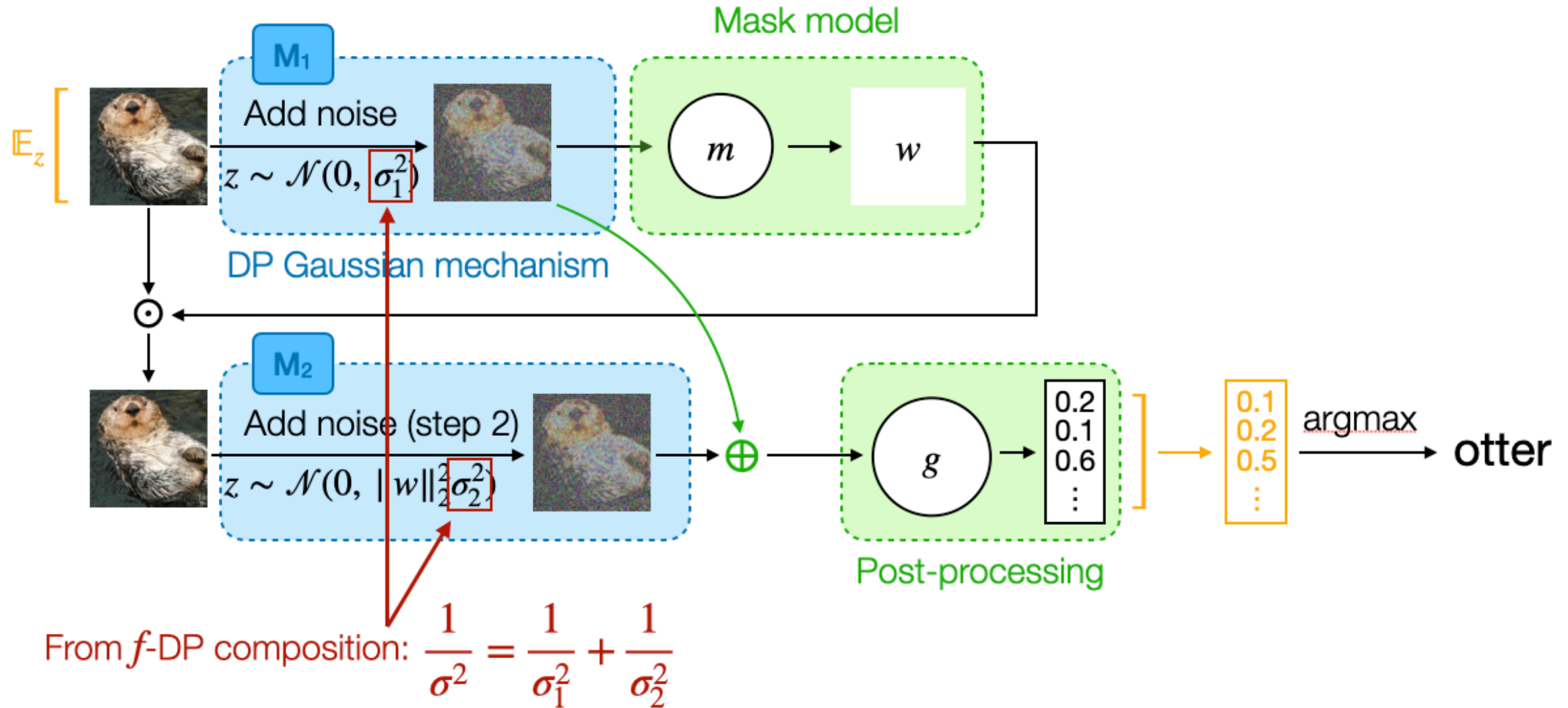
# Adaptive Randomized Smoothing for $\ell_\infty$



$\mathbb{E}_z \left[ \quad \xrightarrow{\substack{\text{Add noise} \\ z \sim \mathcal{N}(0, \sigma_1^2)}} \quad \boxed{M_1} \quad \rightarrow \quad g \rightarrow \begin{matrix} 0.2 \\ 0.1 \\ 0.6 \\ \vdots \end{matrix} \quad \right] \xrightarrow{} \begin{matrix} 0.1 \\ 0.2 \\ 0.5 \\ \vdots \end{matrix} \xrightarrow{\text{argmax}} \text{otter}$

DP Gaussian mechanism  Post-processing

# Adaptive Randomized Smoothing for $\ell_\infty$

# Adaptive Randomized Smoothing for $\ell_\infty$



Mask model

$M_1$

Add noise

$z \sim \mathcal{N}(0, \sigma_1^2)$

DP Gaussian mechanism

$m$ → $w$

$\mathbb{E}_z$

$M_2$

Add noise (step 2)

$z \sim \mathcal{N}(0, \|w\|_2^2 \sigma_2^2)$

$\oplus$

Post-processing

$g$

0.2
0.1
0.6
⋮

0.1
0.2
0.5
⋮

argmax → otter

From $f$-DP composition: $\dfrac{1}{\sigma^2} = \dfrac{1}{\sigma_1^2} + \dfrac{1}{\sigma_2^2}$

# Why does ARS help?



Mask model

$M_1$

Add noise

$z \sim \mathcal{N}(0, \sigma_1^2)$

DP Gaussian mechanism

$m$

$\mathbb{E}_z$

$\odot$

$M_2$

Add noise (step 2)

$z \sim \mathcal{N}(0, \|w\|_2^2 \sigma_2^2)$

$g$

| 0.2 |
| 0.1 |
| 0.6 |
| : |

Post-processing

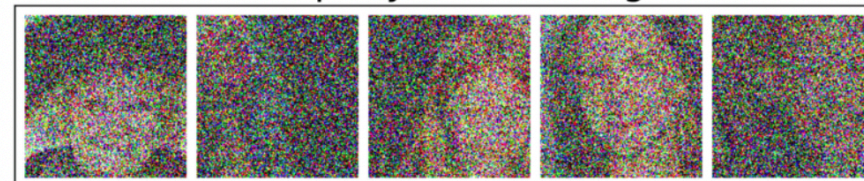| 0.1 |
| 0.2 |
| 0.5 |
| : |

argmax

mouth slightly open

Noise reduction from masking based dimension reduction!
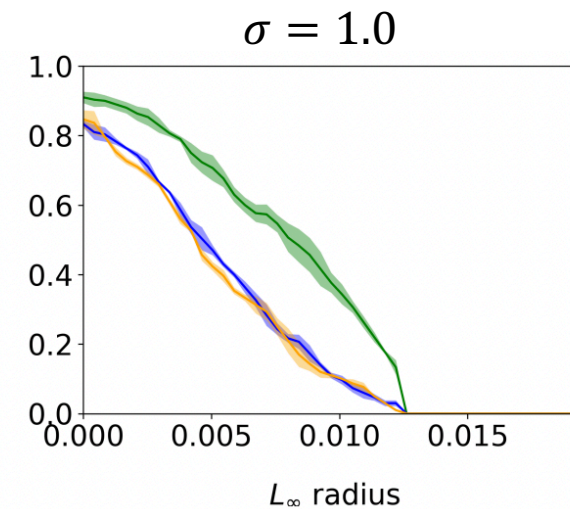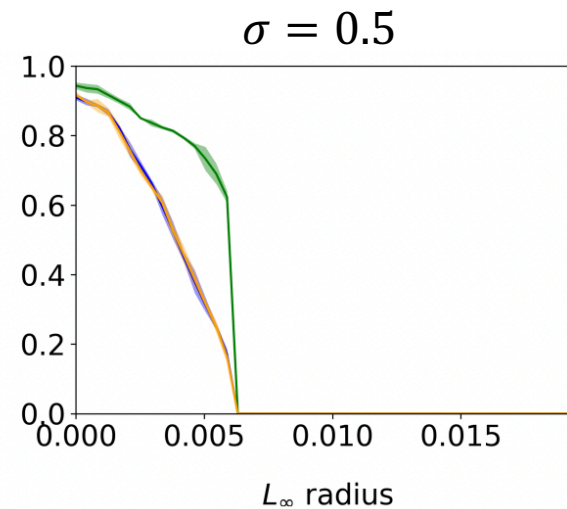
# Evaluations

- **CelebA**



Input images

First query noised images

ARS Masks

Second query noised images after averaging

$\sigma = 0.25$     $\sigma = 0.5$     $\sigma = 1.0$

- ARS (ours)
- Cohen et al.
- static mask

Certified Accuracy vs $L_\infty$ radius

# Conclusion

- Adaptive Randomized Smoothing (ARS) uses DP composition post-processing properties to certify <span style="color:red">adaptive multi-step models</span>.

- ARS learns to <span style="color:red">adjust the scale of noise</span> based on the test input.

- ARS provides <span style="color:red">higher accuracy</span> at a given level of provable robustness.

Link to our code