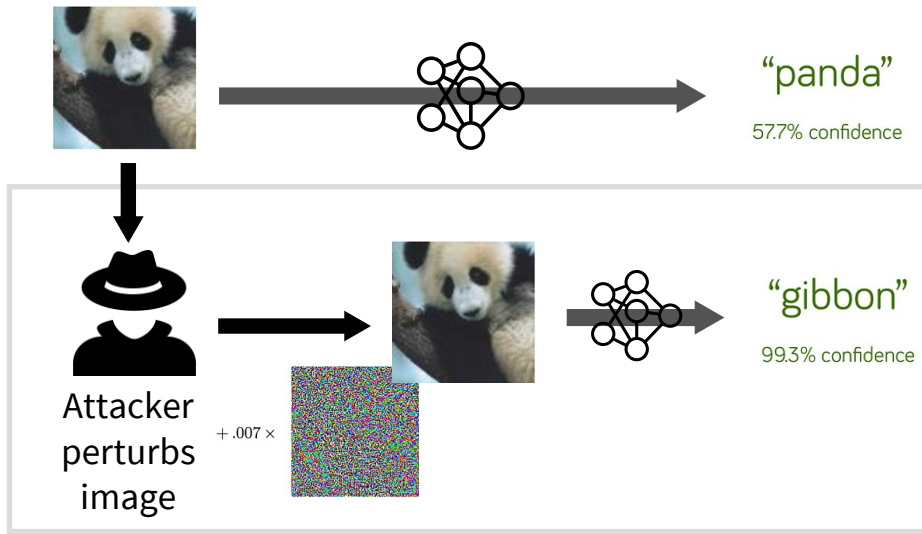# AdaptDel: Adaptable Deletion Rate Randomized Smoothing for Certified Robustness

Zhuoqun Huang, Neil G. Marchant
Olga Ohrimenko, Benjamin I.P. Rubinstein
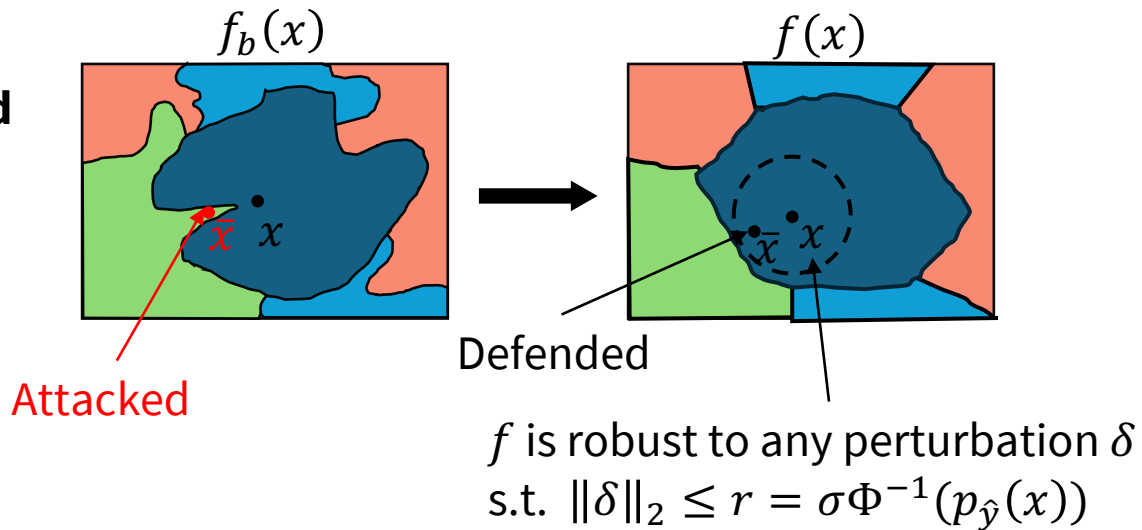
The University of Melbourne

**Attack**



"panda"

57.7% confidence

Attacker perturbs image

$+ .007 \times$

"gibbon"

99.3% confidence

**Certified Defense**

Gaussian Mech. Smoothing (Cohen et al. 2019)

$$f(x) := \arg\max_{y \in Y} \mathbb{E}_{z \sim \mathcal{N}(x, \sigma^2 I)} \left[ \mathbf{1}_{f_\mathbf{b}(z) = y} \right]$$

$f_b(x)$

$f(x)$

Attacked

Defended

$f$ is robust to any perturbation $\delta$
s.t. $\|\delta\|_2 \leq r = \sigma \Phi^{-1}(p_{\hat{y}}(x))$

**Attack**



"panda"
57.7% confidence

Attacker perturbs image
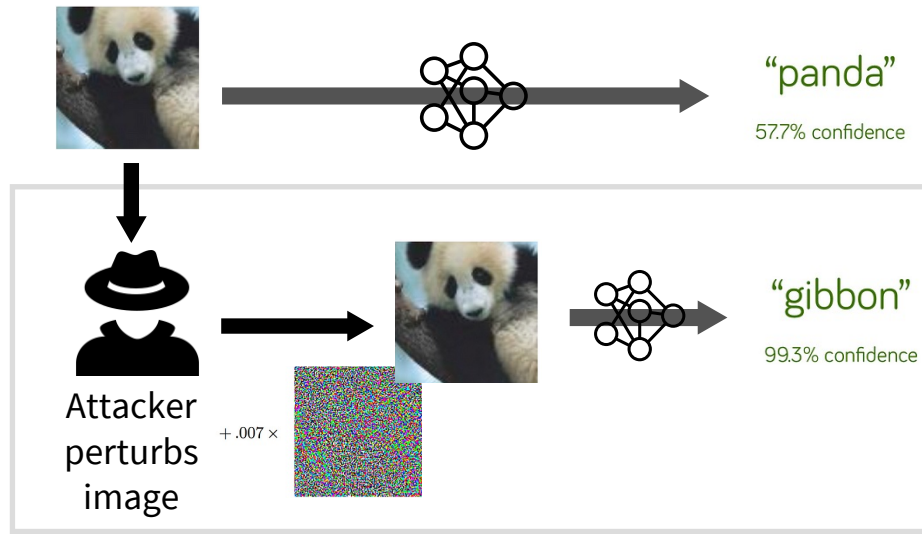$+ .007 \times$

"gibbon"
99.3% confidence

**Certified Defense**

Gaussian Mech. Smoothing (Cohen et al. 2019)

$$f(x) := \arg\max_{y \in Y} \mathbb{E}_{z \sim \mathcal{N}(x,\sigma^2 I)} [\mathbf{1}_{f_b(z) = y}]$$

$f_b(x)$  $f(x)$

Attacked

Defended

$f$ is robust to any perturbation $\delta$
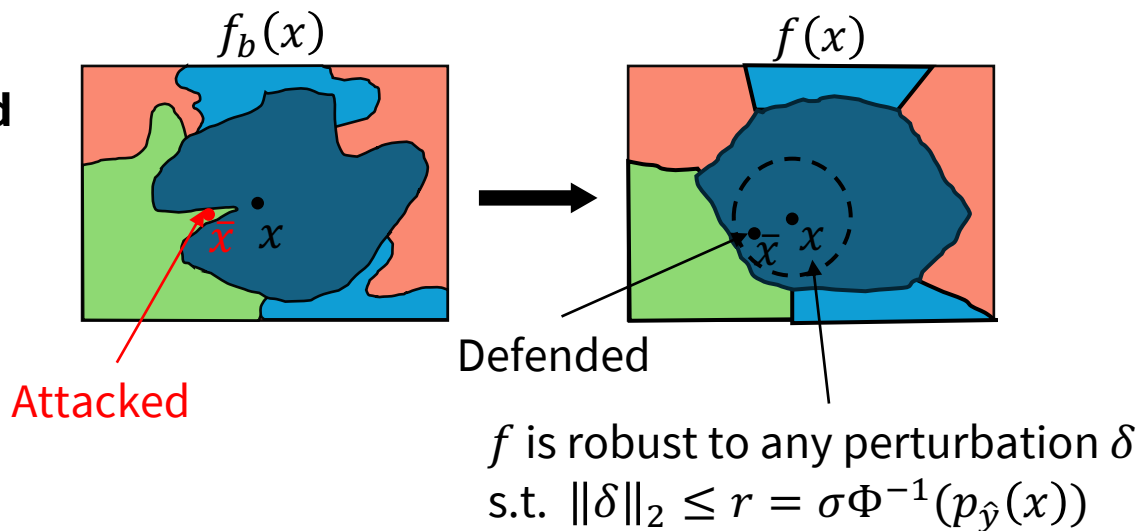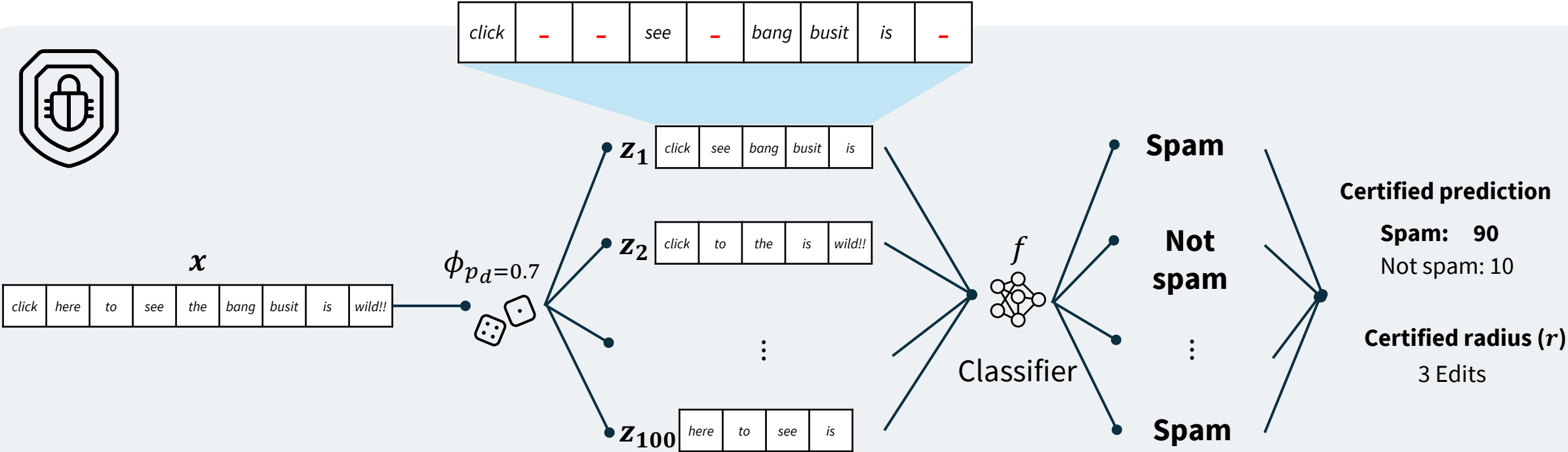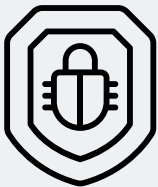s.t. $\|\delta\|_2 \leq r = \sigma \Phi^{-1}(p_{\hat{y}}(x))$

- # Character-level attack
  - *South Africa's historic Soweto township marks its 100th birthday on Tuesday in a (mood) of optimism.* **[World]**
  - *South Africa's historic Soweto township marks its 100th birthday on Tuesday in a (mooP) of optimism.* **[Sci/Tech]**
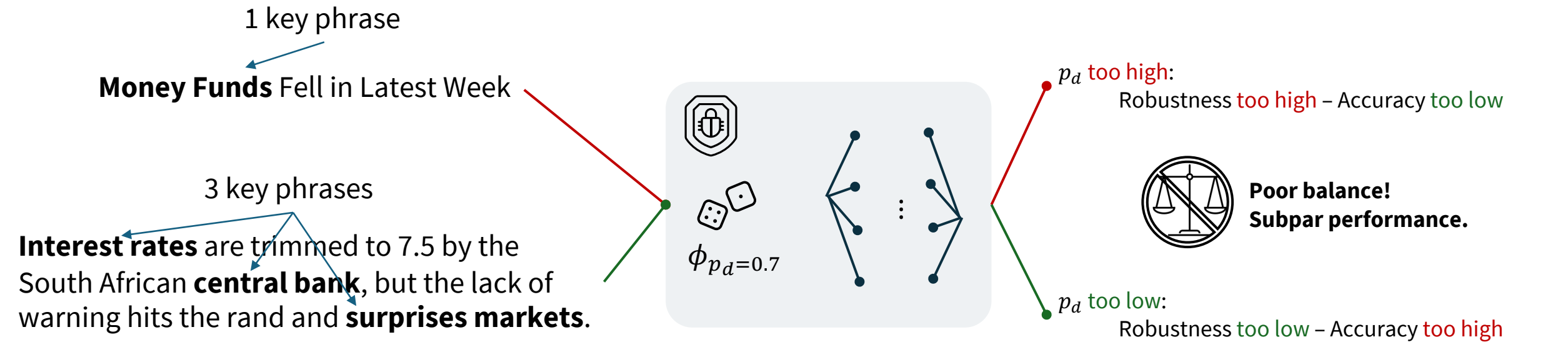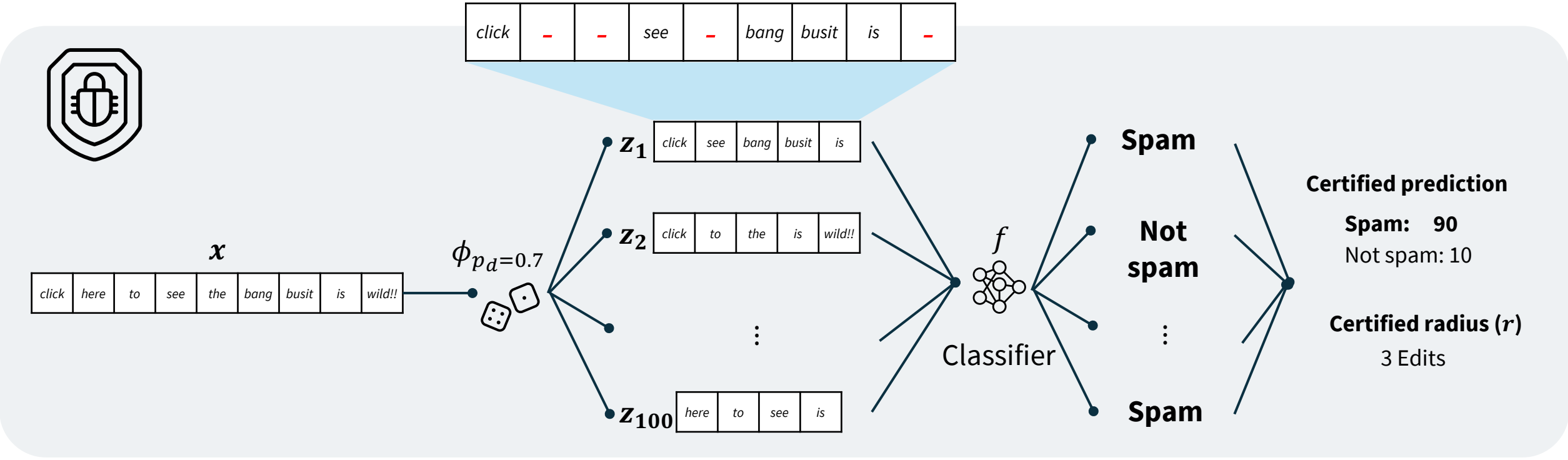    (Ebrahimi et al., ACL 2018)

- # Word-level attack
  - *Super ant colony hits Australia. A (giant) 100km colony of ants could threaten (local) insect species.* **[World]**
  - *Super ant colony hits Australia (Coast). A (gigantic) 100km colony of ants could threaten (local) insect species.* **[Sci/Tech]**
    (Li et al., NAACL 2021)

# Background: CERT-ED
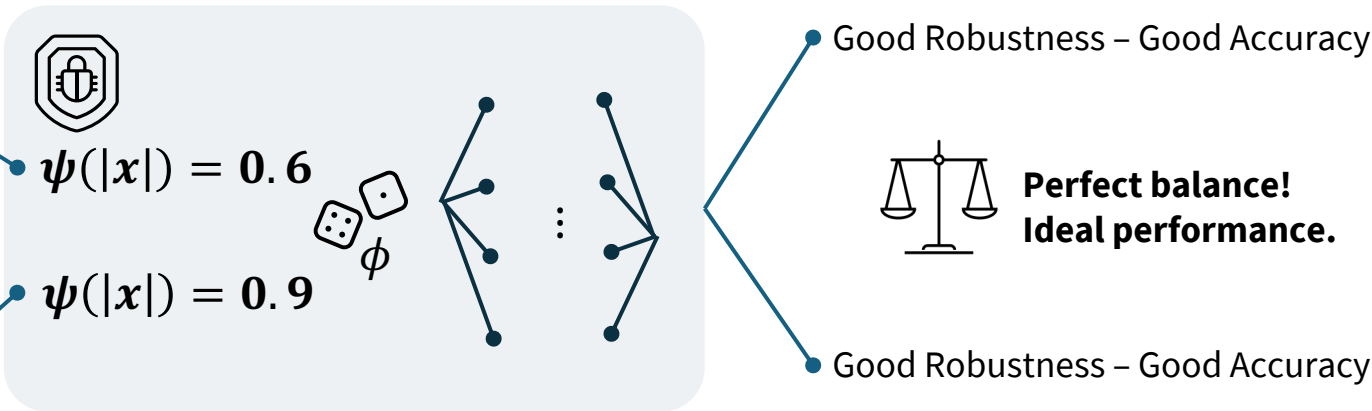
(Huang et al., EMNLP 2024)

# Background: CERT-ED

| click | – | – | see | – | bang | busit | is | – |

$z_1$ | click | see | bang | busit | is |

$z_2$ | click | to | the | is | wild!! |

$x$

| click | here | to | see | the | bang | busit | is | wild!! |

$\phi_{p_d=0.7}$

$z_{100}$ | here | to | see | is |

$f$

Classifier

**Spam**

**Not spam**

**Spam**

**Certified prediction**

**Spam:** 90
Not spam: 10

**Certified radius ($r$)**
3 Edits

---

1 key phrase

**Money Funds** Fell in Latest Week

3 key phrases

**Interest rates** are trimmed to 7.5 by the South African **central bank**, but the lack of warning hits the rand and **surprises markets**.

$\phi_{p_d=0.7}$

$p_d$ too high:
Robustness too high – Accuracy too low

**Poor balance!
Subpar performance.**

$p_d$ too low:
Robustness too low – Accuracy too high

# Model deletion rates as a function of the input sequence $p_{del} = \psi(|x|)$

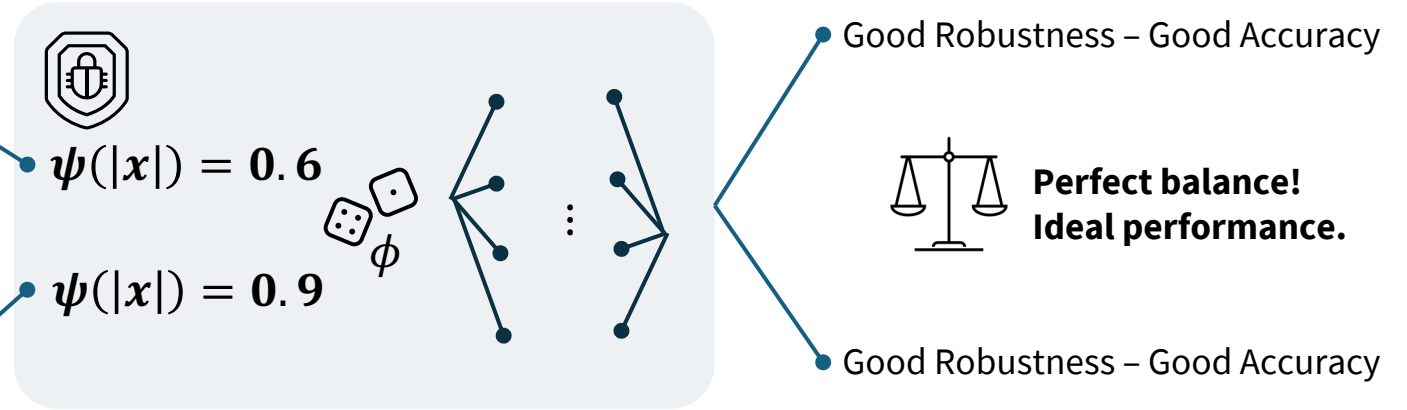Money Funds Fell in Latest Week

$\psi(|x|) = 0.6$

$\psi(|x|) = 0.9$

$\phi$

Interest rates are trimmed to 7.5 by the South African central bank, but the lack of warning hits the rand and surprises markets.

Good Robustness – Good Accuracy

**Perfect balance! Ideal performance.**

Good Robustness – Good Accuracy

# Model deletion rates as a function of the input sequence $p_{del} = \psi(|x|)$

Money Funds Fell in Latest Week

$\psi(|x|) = 0.6$

$\phi$

Good Robustness – Good Accuracy

**Perfect balance!
Ideal performance.**

Interest rates are trimmed to 7.5 by the South African central bank, but the lack of warning hits the rand and surprises markets.

$\psi(|x|) = 0.9$

Good Robustness – Good Accuracy

---

**Algorithm 1** CERTIFY

**Require:** base classifier $f_{\mathrm{b}}$, input sequence $x$, predicted class $y_1$, length-dependent deletion probability $\psi$, allowed edit operations o, significance level $\alpha$

**Ensure:** maximum radius that can be certified

1: $t_1^{\mathrm{lb}} \leftarrow \hat{p}_{y_1}^{\mathrm{lb}}(x; f_{\mathrm{b}}, \phi_\psi, \alpha)$
2: $t_2^{\mathrm{ub}} \leftarrow \max_{y \neq y_1} \hat{p}_y^{\mathrm{ub}}(x; f_{\mathrm{b}}, \phi_\psi, \alpha)$
3: **for** $r = 0$ **to** $\infty$ **do**
4:     **for all** $(n_{\mathrm{del}}, n_{\mathrm{ins}}, n_{\mathrm{sub}}) \in \mathcal{C}(\mathrm{o}, r)$ **do**
5:         $|\bar{x}| \leftarrow |x| + n_{\mathrm{ins}} - n_{\mathrm{del}}$
6:         $\bar{t}_1^{\mathrm{lb}} \leftarrow \mathrm{lb}(t_1^{\mathrm{lb}}, x, \bar{x}, \psi)$
7:         $\bar{t}_2^{\mathrm{ub}} \leftarrow \mathrm{ub}(t_2^{\mathrm{ub}}, x, \bar{x}, \psi)$
8:         **if** $\bar{t}_1^{\mathrm{lb}} \leq \bar{t}_2^{\mathrm{ub}}$ **then**
9:             **return** $r$

# Model deletion rates as a function of the input sequence $p_{del} = \psi(|x|)$

Money Funds Fell in Latest Week

Interest rates are trimmed to 7.5 by the South African central bank, but the lack of warning hits the rand and surprises markets.
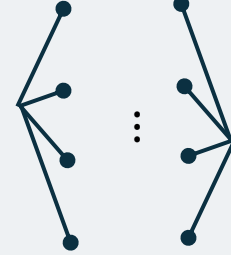
$\psi(|x|) = 0.6$

$\psi(|x|) = 0.9$

$\phi$

Good Robustness – Good Accuracy

**Perfect balance! Ideal performance.**

Good Robustness – Good Accuracy

---

**Algorithm 1** CERTIFY

**Require:** base classifier $f_b$, input sequence $x$, predicted class $y_1$, length-dependent deletion probability $\psi$, allowed edit operations o, significance level $\alpha$

**Ensure:** maximum radius that can be certified
1: $t_1^{lb} \leftarrow \hat{p}_{y_1}^{lb}(x; f_b, \phi_\psi, \alpha)$
2: $t_2^{ub} \leftarrow \max_{y \neq y_1} \hat{p}_y^{ub}(x; f_b, \phi_\psi, \alpha)$
3: **for** $r = 0$ **to** $\infty$ **do**
4:     **for all** $(n_{del}, n_{ins}, n_{sub}) \in \mathcal{C}(o, r)$ **do**
5:        $|\bar{x}| \leftarrow |x| + n_{ins} - n_{del}$
6:        $\bar{t}_1^{lb} \leftarrow lb(t_1^{lb}, x, \bar{x}, \psi)$
7:        $\bar{t}_2^{ub} \leftarrow ub(t_2^{ub}, x, \bar{x}, \psi)$
8:        **if** $\bar{t}_1^{lb} \leq \bar{t}_2^{ub}$ **then**
9:           **return** $r$

**Lemma 3.** *Let* $x, \bar{x} \in \mathcal{X}$ *be a pair of inputs with a longest common subsequence (LCS)* $z^\star$ *and let* $\mu = p_y(x; f_b)$. *Define*

$$H^* = \begin{cases} \min_{h:\sum_{i=0}^{h} \mathcal{B}_i(|z^\star|,\psi) \geq \mu-1+\psi^{|x|-|z^\star|}} h, & \psi \geq \bar{\psi}, \\ \max_{h:\sum_{i=h}^{|z^\star|} \mathcal{B}_i(|z^\star|,\psi) \geq \mu-1+\psi^{|x|-|z^\star|}} h, & \psi < \bar{\psi}, \end{cases}$$

*as a threshold on the number of tokens retained when editing* $x$, *where* $\mathcal{B}_k(n, p) := \binom{n}{k}(1-p)^k p^{n-k}$ *is the Binomial pmf for* $n$ *trials with success probability* $1 - p$. *Then there exists a lower bound* $lb(\mu, x, \bar{x}, \psi) \leq p_y(\bar{x}; f_b)$ *such that:*

$$lb(\mu, x, \bar{x}, \psi) = \frac{\bar{\psi}^{|\bar{x}|-|z^\star|}}{\psi^{|x|-|z^\star|}} \left( \sum_{i=l(H^*+1)}^{(1-l)(H^*-1)+l|z^\star|} \mathcal{B}_i(|z^\star|, \bar{\psi}) \right.$$
$$\left. + \mathcal{B}_{H^*}(|z^\star|, \bar{\psi}) \left\lfloor \frac{c(\mu, |x|, |z^\star|, \psi, H^*)}{\mathcal{B}_{H^*}(|z^\star|, \psi)} \right\rfloor_{\binom{|z^\star|}{H^*}^{-1}} \right),$$

*where*

$$c(\mu, |x|, |z^\star|, \psi, H^*) = \mu - 1 + \psi^{|x|-|z^\star|} - \sum_{i=l(H^*+1)}^{(1-l)(H^*-1)+l|z^\star|} \mathcal{B}_j(|z^\star|, \psi),$$

$l = \mathbf{1}_{\psi < \bar{\psi}}$ *is a binary indicator and* $\lfloor \cdot \rfloor_v := \lfloor \frac{\cdot}{v} \rfloor v$ *is a gridded flooring operation.*
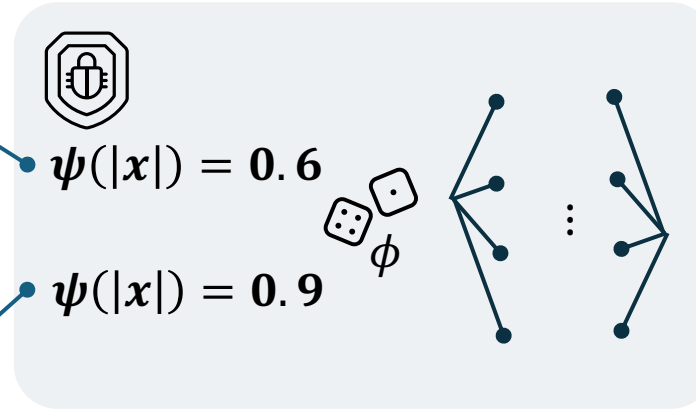
(Huang et al., NeurIPS 2023)

Sketch: Start with LCS technique, and progressively lower bound the quantity. Finally, formulate the minimization problem as knapsack problem.

# Model deletion rates as a function of the input sequence $p_{del} = \psi(|x|)$

Money Funds Fell in Latest Week

$$\psi(|x|) = 0.6$$

$$\psi(|x|) = 0.9$$

$\phi$

Interest rates are trimmed to 7.5 by the South African central bank, but the lack of warning hits the rand and surprises markets.

Good Robustness – Good Accuracy

**Perfect balance! Ideal performance.**

Good Robustness – Good Accuracy

---

**Algorithm 1** CERTIFY

**Require:** base classifier $f_b$, input sequence $x$, predicted class $y_1$, length-dependent deletion probability $\psi$, allowed edit operations o, significance level $\alpha$

**Ensure:** maximum radius that can be certified
1: $t_1^{lb} \leftarrow \hat{p}_{y_1}^{lb}(x; f_b, \phi_\psi, \alpha)$
2: $t_2^{ub} \leftarrow \max_{y \neq y_1} \hat{p}_y^{ub}(x; f_b, \phi_\psi, \alpha)$
3: **for** $r = 0$ **to** $\infty$ **do**
4:     **for all** $(n_{del}, n_{ins}, n_{sub}) \in \mathcal{C}(o, r)$ **do**
5:       $|\bar{x}| \leftarrow |x| + n_{ins} - n_{del}$
6:       $\bar{t}_1^{lb} \leftarrow lb(t_1^{lb}, x, \bar{x}, \psi)$
7:       $\bar{t}_2^{ub} \leftarrow ub(t_2^{ub}, x, \bar{x}, \psi)$
8:       **if** $\bar{t}_1^{lb} \leq \bar{t}_2^{ub}$ **then**
9:         **return** $r$

**Lemma 3.** *Let* $x, \bar{x} \in \mathcal{X}$ *be a pair of inputs with a longest common subsequence (LCS)* $z^\star$ *and let* $\mu = p_y(x; f_b)$. *Define*

$$H^* = \begin{cases} \min_{h: \sum_{i=0}^h \mathcal{B}_i(|z^\star|, \psi) \geq \mu - 1 + \psi^{|x| - |z^\star|}} h, & \psi \geq \bar{\psi}, \\ \max_{h: \sum_{i=h}^{|z^\star|} \mathcal{B}_i(|z^\star|, \psi) \geq \mu - 1 + \psi^{|x| - |z^\star|}} h, & \psi < \bar{\psi}, \end{cases}$$

*as a threshold on the number of tokens retained when editing* $x$, *where* $\mathcal{B}_k(n, p) := \binom{n}{k}(1-p)^k p^{n-k}$ *is the Binomial pmf for* $n$ *trials with success probability* $1 - p$. *Then there exists a lower bound* $lb(\mu, x, \bar{x}, \psi) \leq p_y(\bar{x}; f_b)$ *such that:*

$$lb(\mu, x, \bar{x}, \psi) = \frac{\bar{\psi}^{|\bar{x}| - |z^\star|}}{\psi^{|x| - |z^\star|}} \left( \sum_{i=l(H^*+1)}^{(1-l)(H^*-1)+l|z^\star|} \mathcal{B}_i(|z^\star|, \bar{\psi}) \right.$$

$$\left. + \mathcal{B}_{H^*}(|z^\star|, \bar{\psi}) \left\lfloor \frac{c(\mu, |x|, |z^\star|, \psi, H^*)}{\mathcal{B}_{H^*}(|z^\star|, \psi)} \right\rfloor_{\binom{|z^\star|}{H^*}^{-1}} \right),$$

*where*

$$c(\mu, |x|, |z^\star|, \psi, H^*) = \mu - 1 + \psi^{|x| - |z^\star|} - \sum_{i=l(H^*+1)}^{(1-l)(H^*-1)+l|z^\star|} \mathcal{B}_j(|z^\star|, \psi),$$

$l = \mathbf{1}_{\psi < \bar{\psi}}$ *is a binary indicator and* $\lfloor \cdot \rfloor_v := \lfloor \frac{\cdot}{v} \rfloor v$ *is a gridded flooring operation.*

(Huang et al., NeurIPS 2023)

Sketch: Start with LCS technique, and progressively lower bound the quantity. Finally, formulate the maximization problem as knapsack problem.
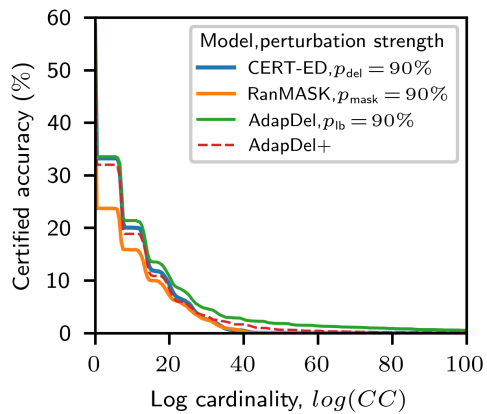
**AdaptDel**

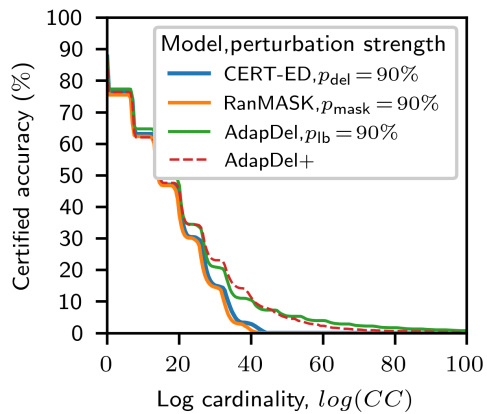$$\psi(|x|) = \max\left(p_{lb}, 1 - \frac{k}{|x|}\right)$$

**AdaptDel+**

$\psi(|x|)$ calibrated automatically by length binning and golden section search
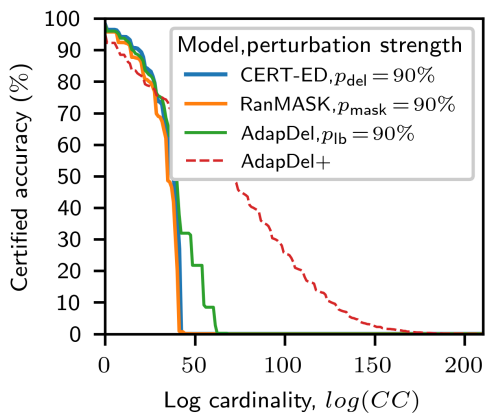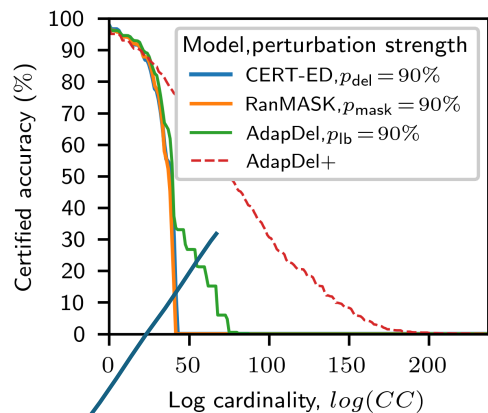
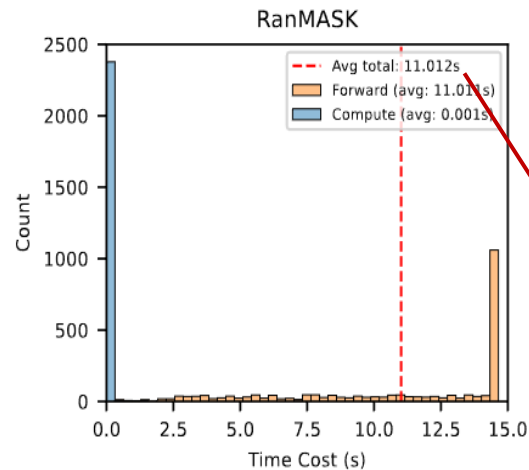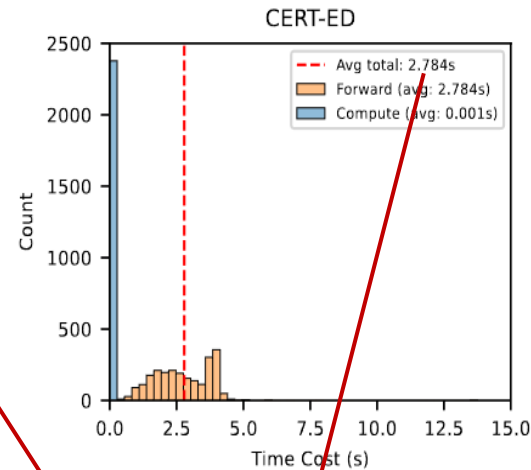Improved certified accuracy on all datasets

(a) Yelp
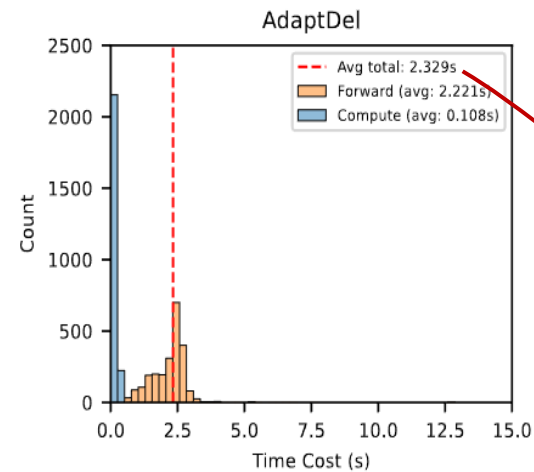
(b) IMDB

(c) LUN

(d) Spam Assassin

Significant improvement adapting to longer sequences

(a) RanMASK

(b) CERT-ED

(c) AdaptDel

20% and 372% faster