



ADVERSARIAL OBSERVATIONS IN WEATHER FORECASTING

Erik Imgrund, Thorsten Eisenhofer, Konrad Rieck

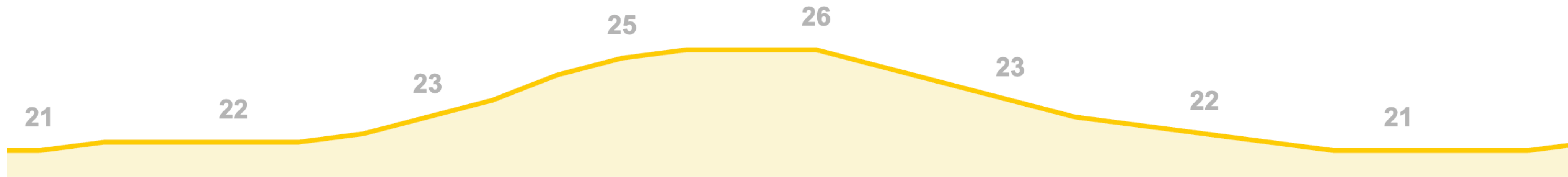


22 °C | °F

Precipitation: 90%
Humidity: 70%
Wind: 17 km/h

Weather
Thursday
Light rain

Temperature | Precipitation | Wind



03:00

06:00

09:00

12:00

15:00

18:00

21:00

00:00

Mon



31° 26°

Tue



31° 26°

Wed



27° 24°

Thu



22° 22°

Fri



26° 23°

Sat



27° 23°

Sun



26° 23°

Mon



26° 23°

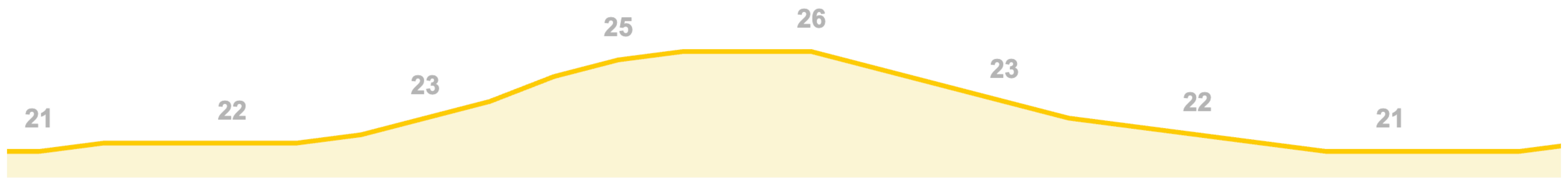


~~22~~ °C | °F

Precipitation: 90%
Humidity: 70%
Wind: 17 km/h

Weather
Thursday
Light rain

Temperature | Precipitation | Wind



03:00

06:00

09:00

12:00

15:00

18:00

21:00

00:00

Mon

Tue

Wed

Thu

Fri

Sat

Sun

Mon



31° 26°



31° 26°



27° 24°



22° 22°



26° 23°



27° 23°



26° 23°



26° 23°

Can we manipulate the weather?

Can we manipulate the weather?

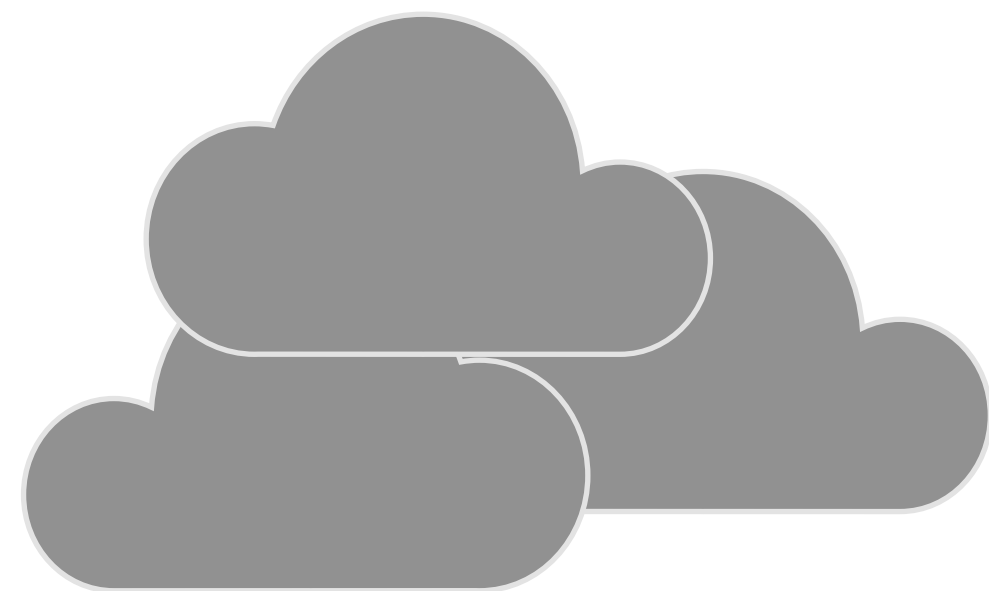
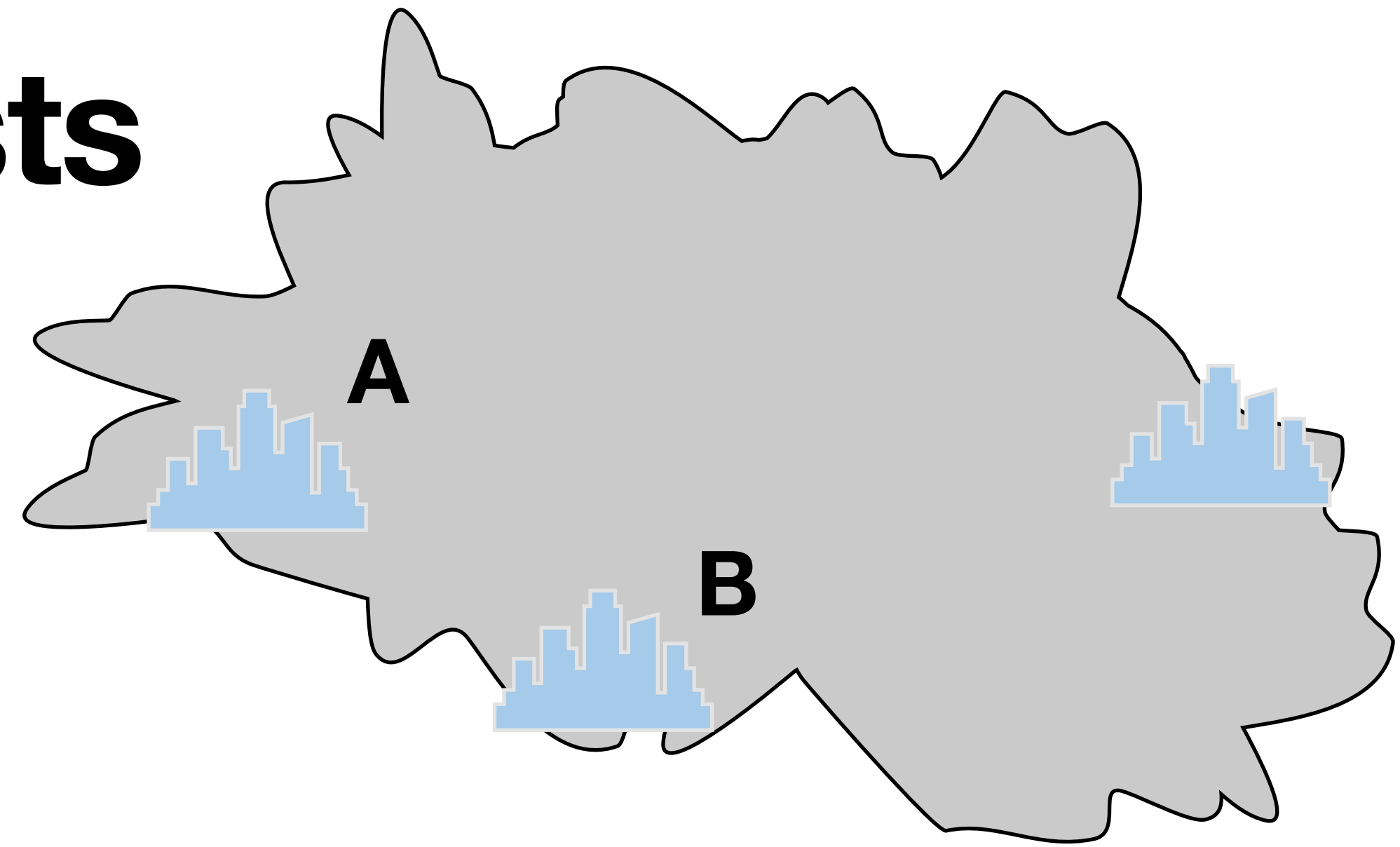
no

forecast

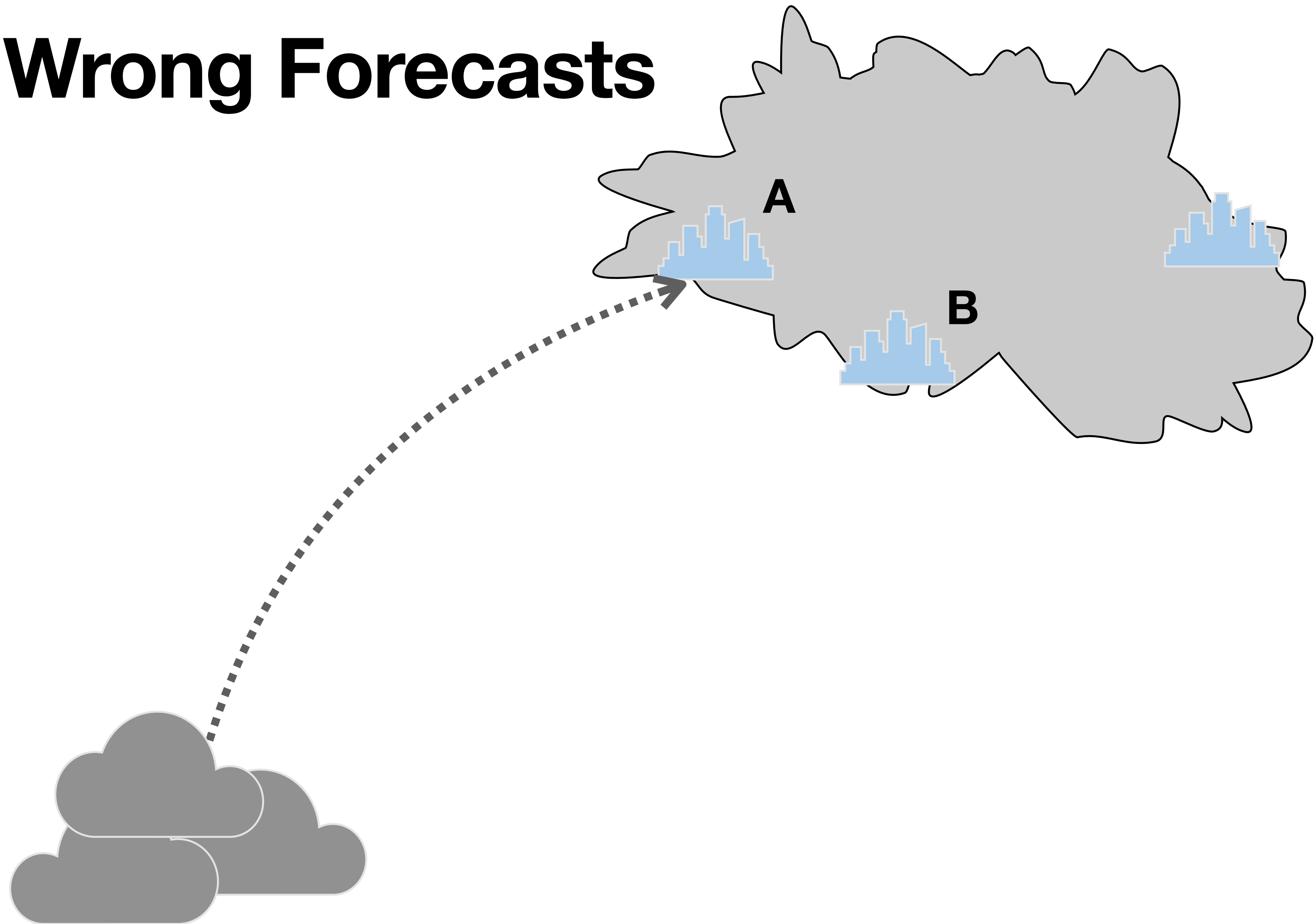
Can we manipulate the weather?

Risk of Wrong Forecasts

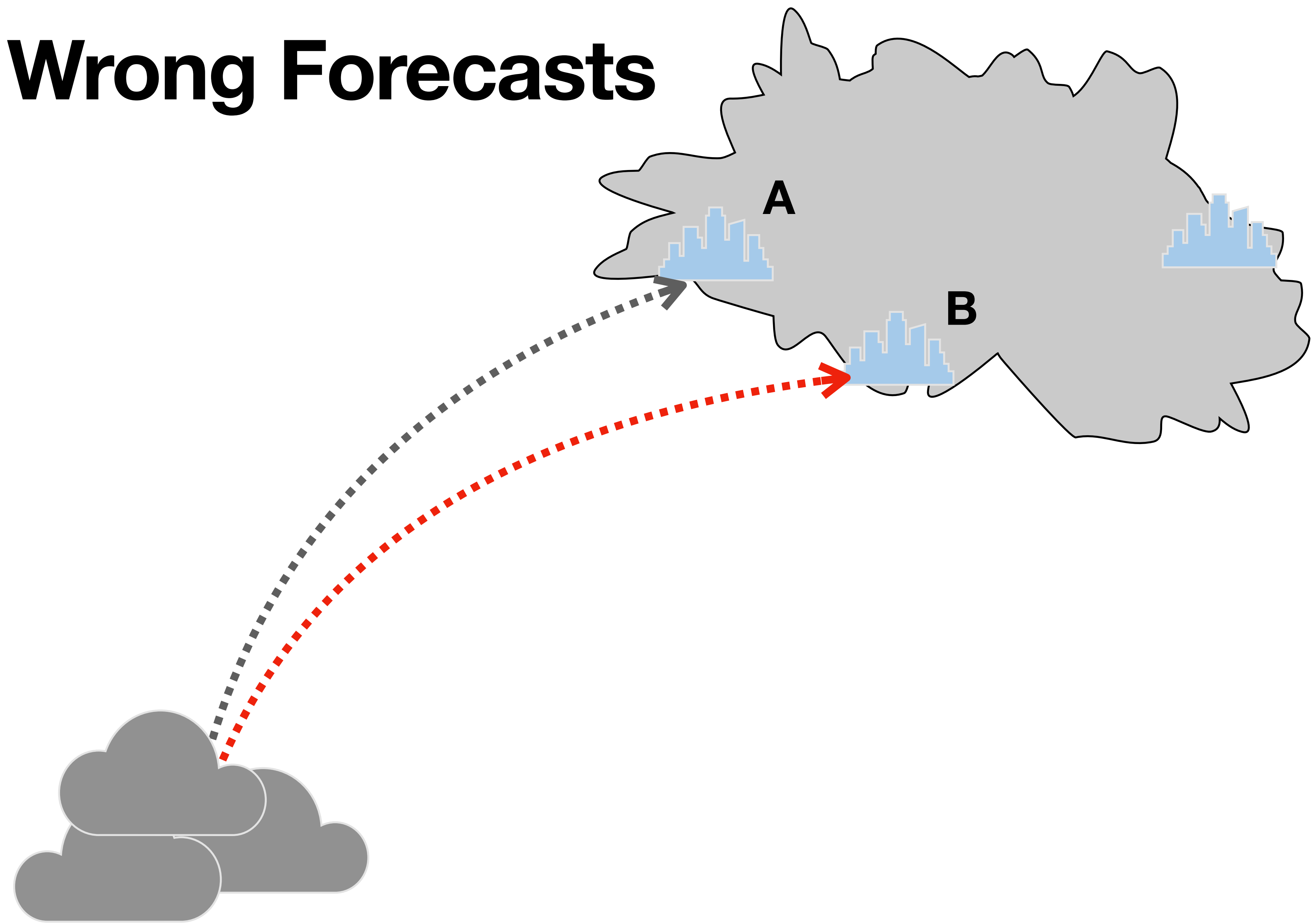
Risk of Wrong Forecasts



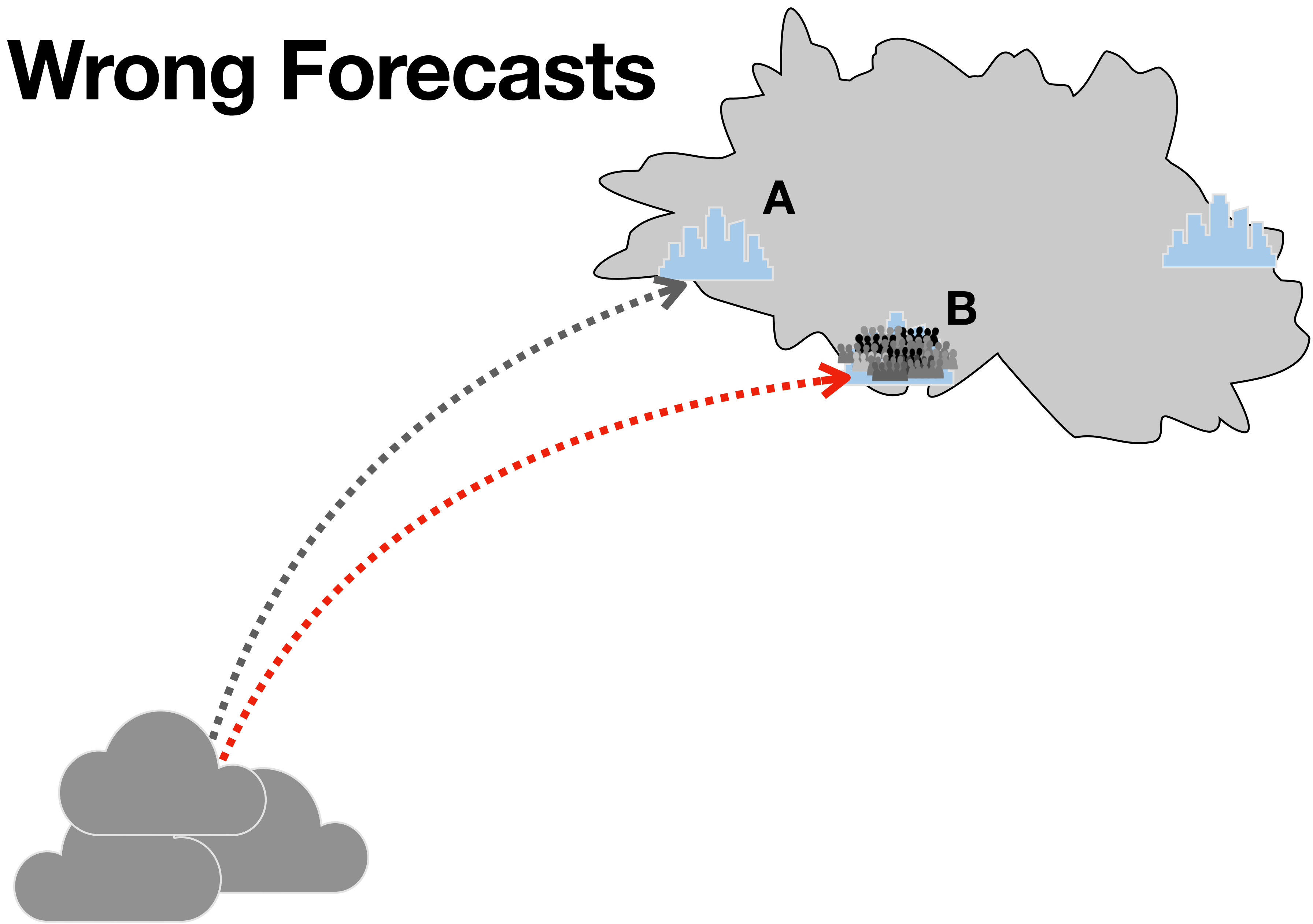
Risk of Wrong Forecasts



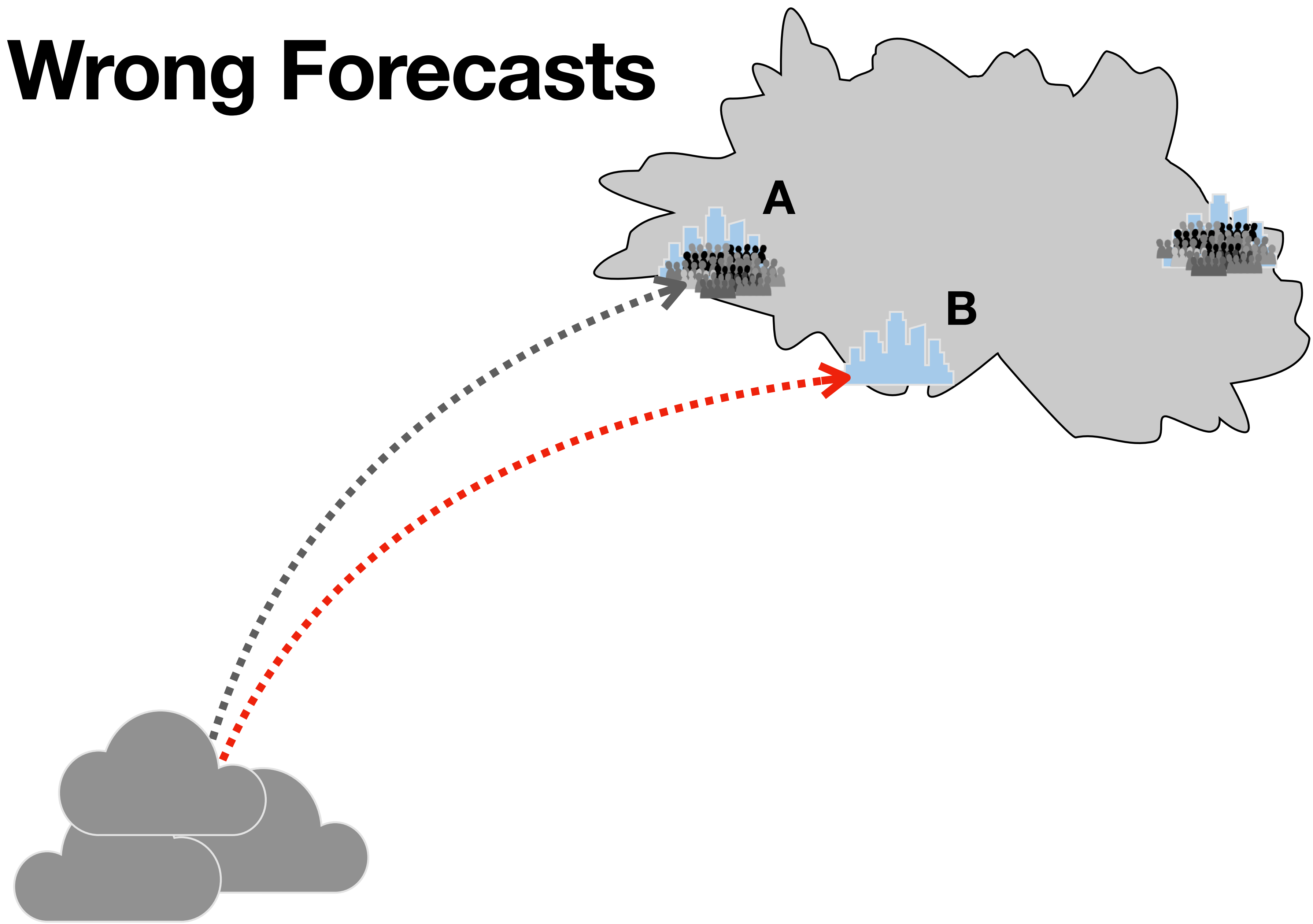
Risk of Wrong Forecasts



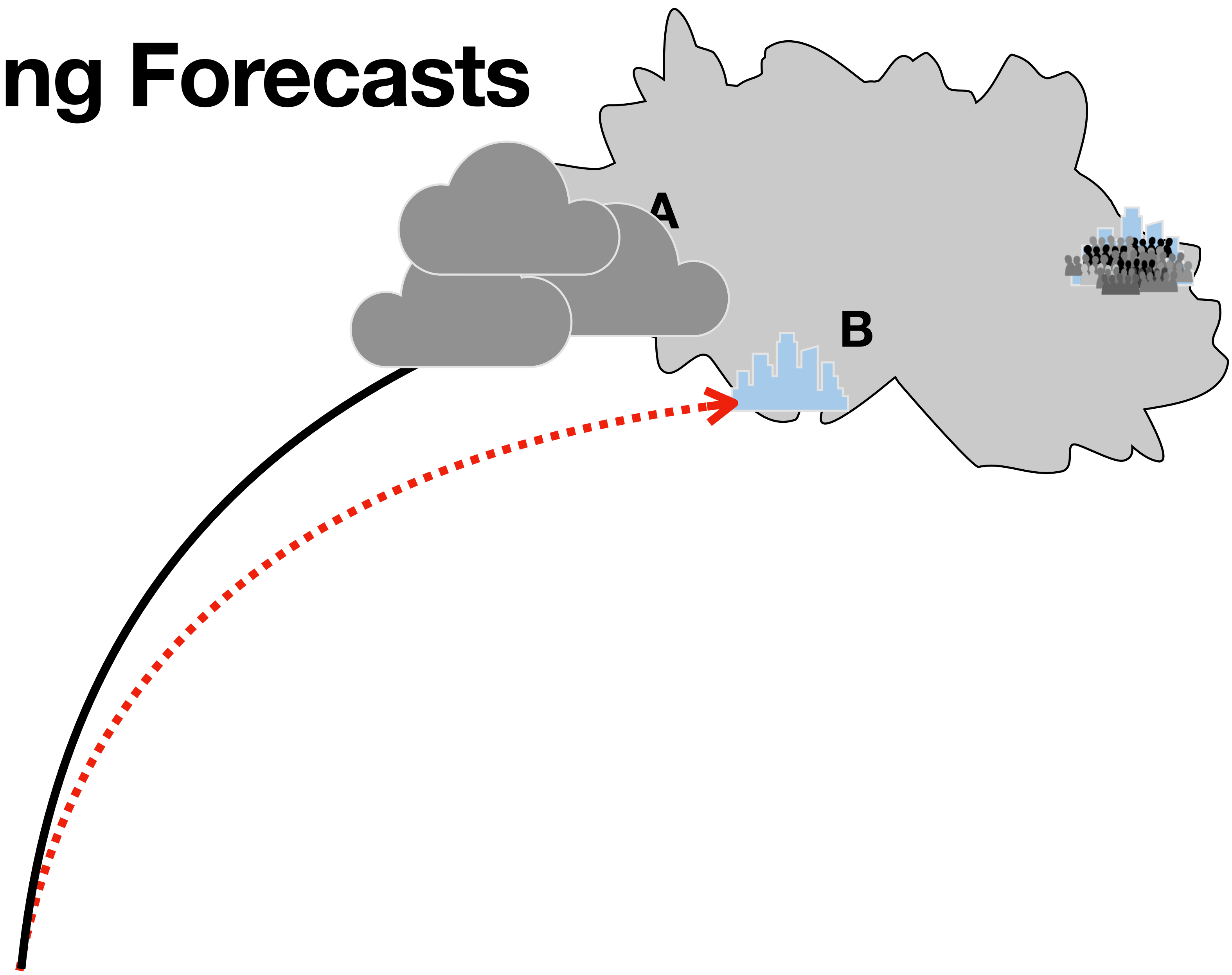
Risk of Wrong Forecasts



Risk of Wrong Forecasts

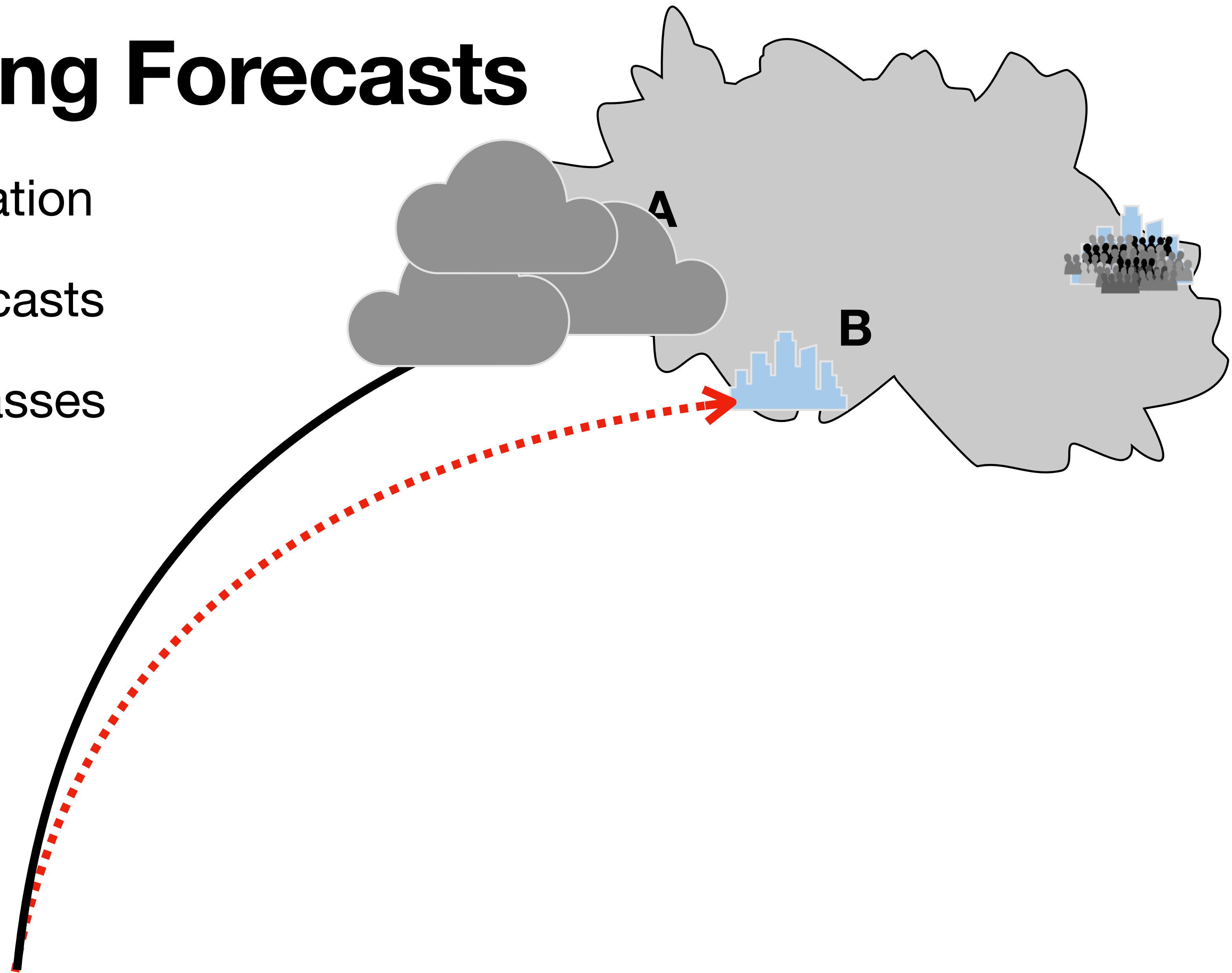


Risk of Wrong Forecasts



Risk of Wrong Forecasts

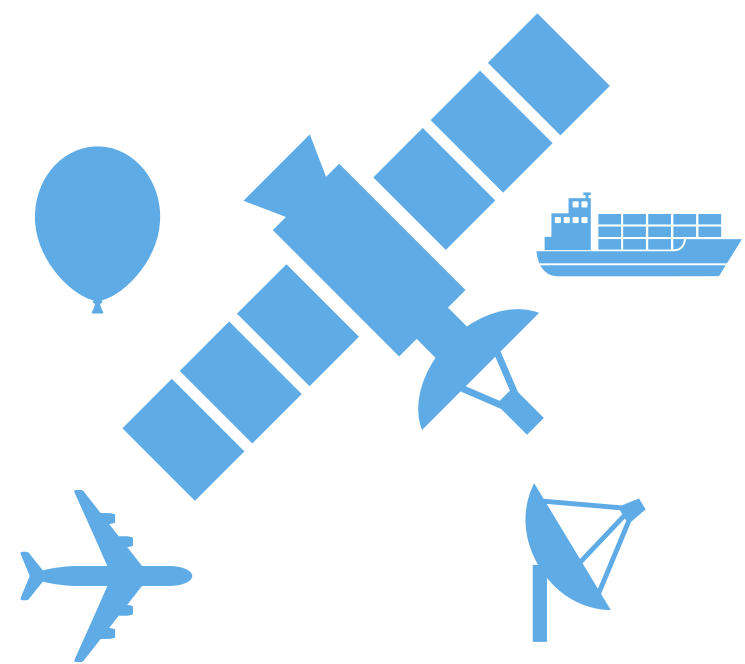
- Can hinder preparation
- Lower trust in forecasts
- Induce panic in masses



Medium-range Weather Forecasting

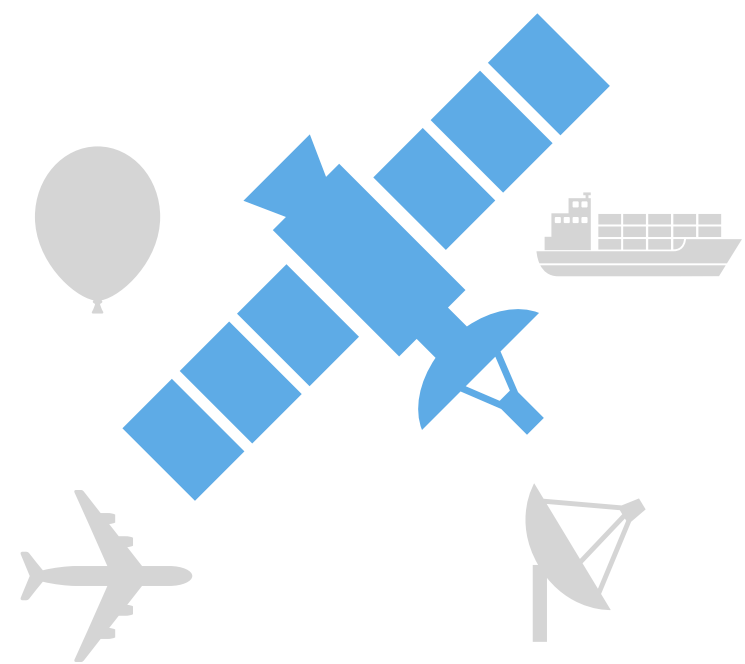
Medium-range Weather Forecasting

- Based on diverse sensor data



Medium-range Weather Forecasting

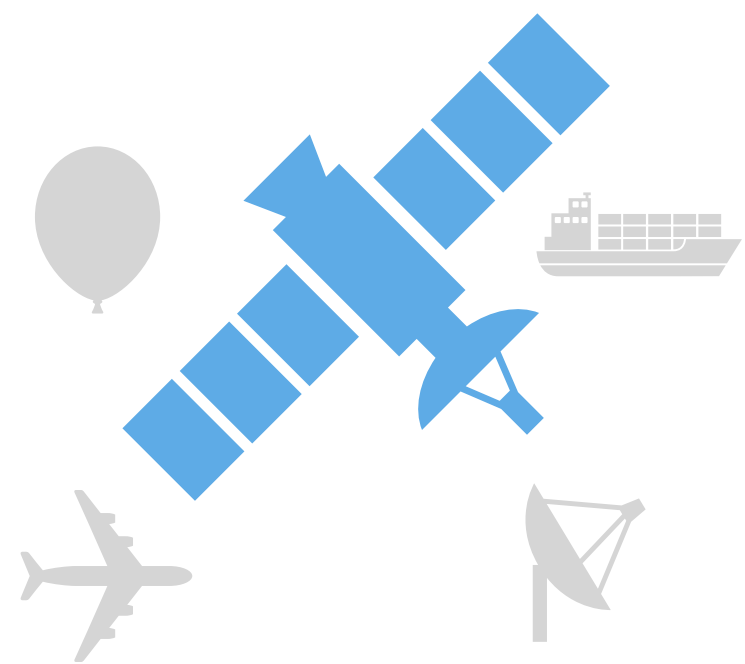
- Based on diverse sensor data
 - Satellites have by far the largest impact¹



¹ ECMWF “Atmospheric Model Data Sources”, Forecast User Guide

Medium-range Weather Forecasting

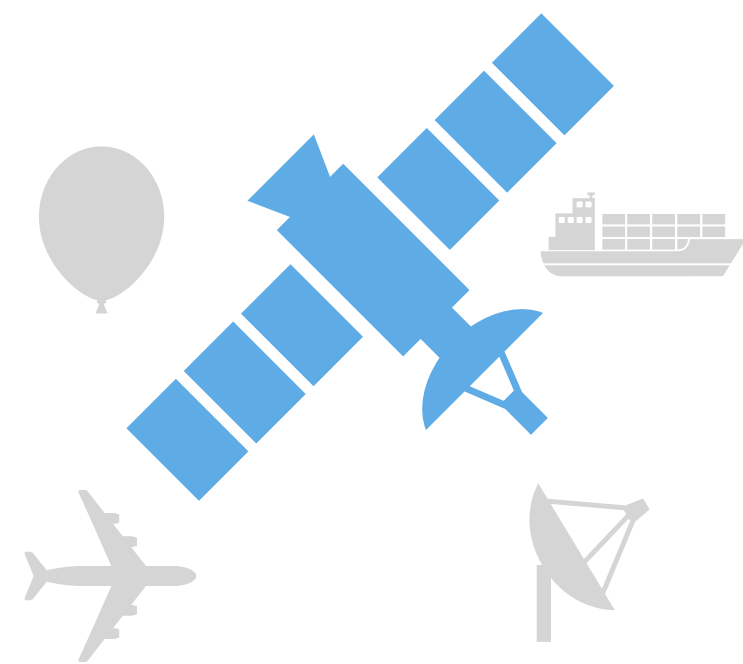
- Based on diverse sensor data
 - Satellites have by far the largest impact¹
- Recently: AI-based forecasting outperforms traditional forecasting



¹ ECMWF “Atmospheric Model Data Sources”, Forecast User Guide

Medium-range Weather Forecasting

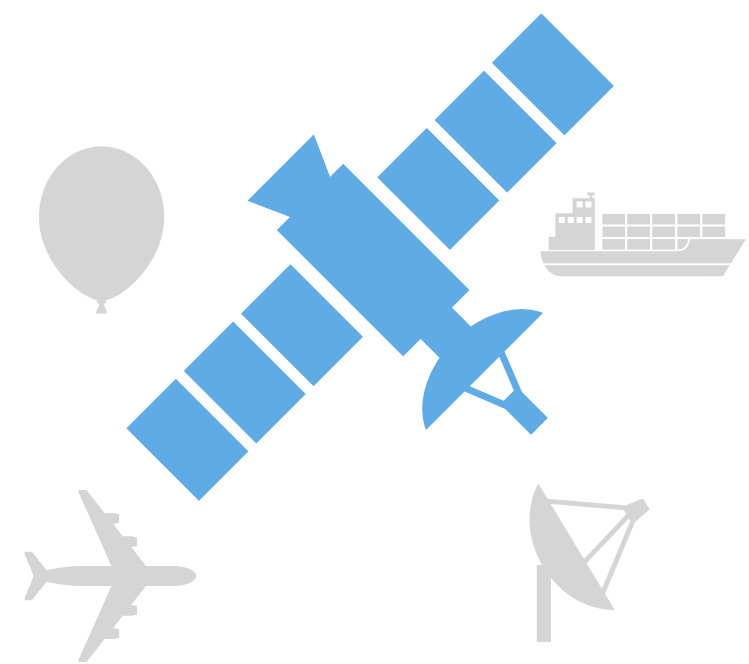
- Based on diverse sensor data
 - Satellites have by far the largest impact¹
- Recently: AI-based forecasting outperforms traditional forecasting
- SotA: Autoregressive diffusion model GenCast²



¹ ECMWF “Atmospheric model data sources”, Forecast User Guide

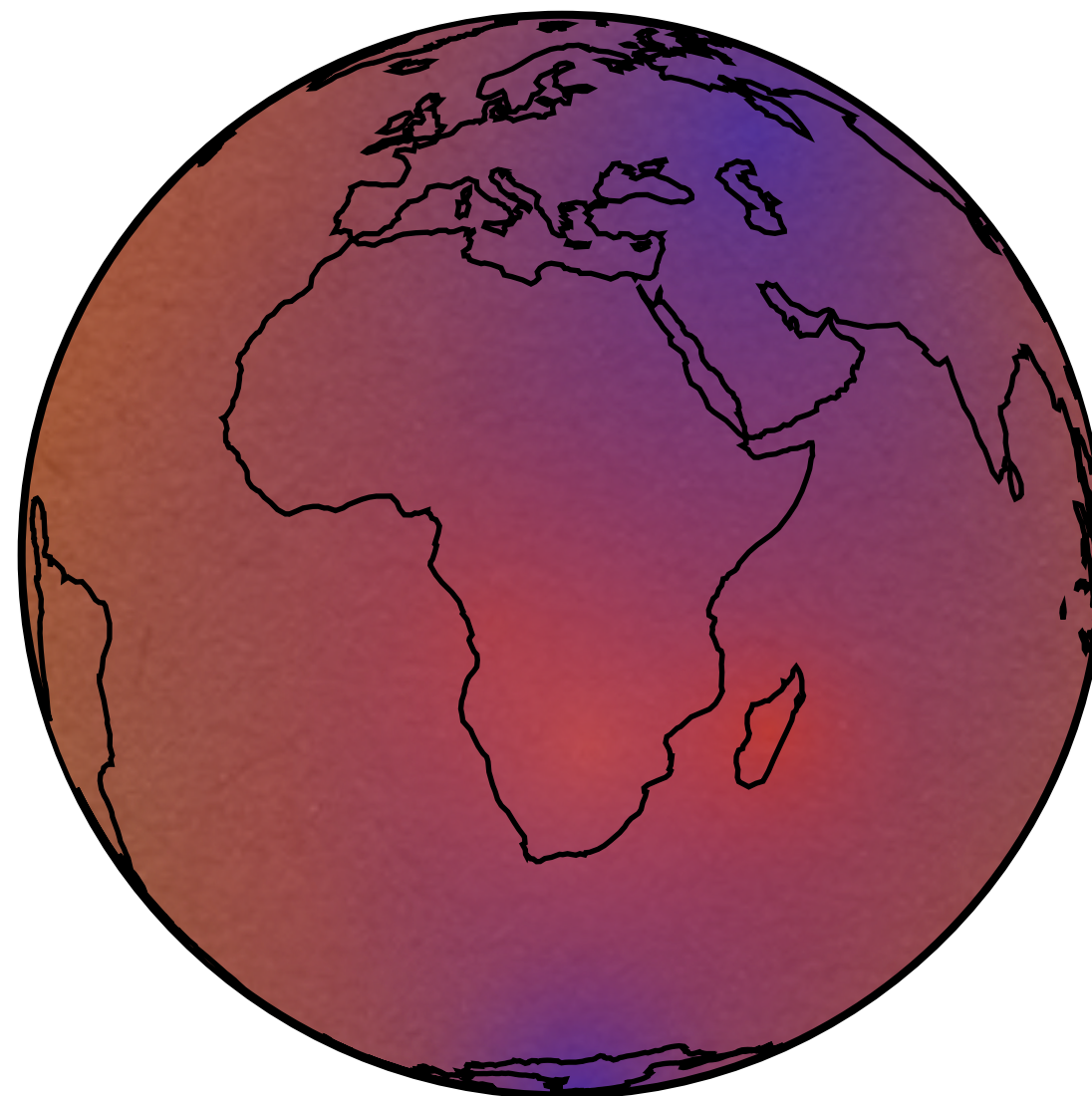
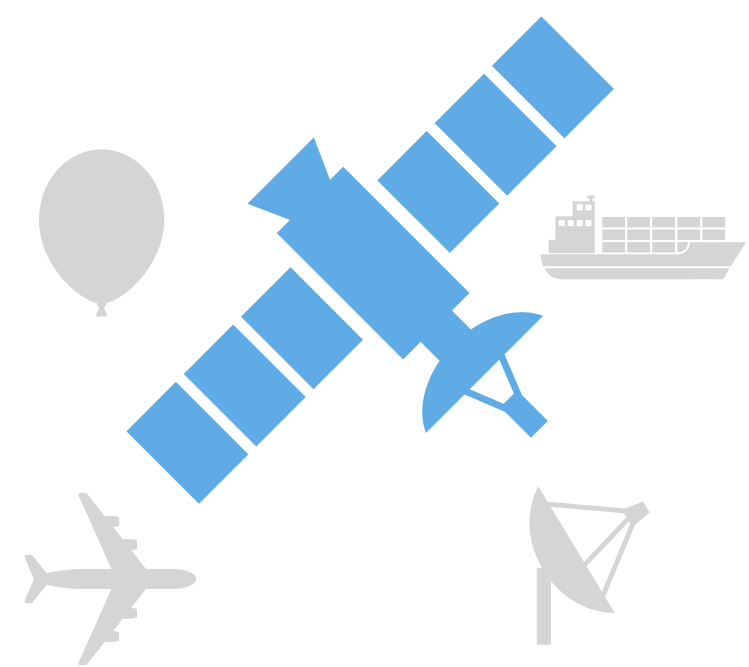
² Price et al. “Probabilistic weather forecasting with machine learning”, Nature

Medium-range Weather Forecasting



Medium-range Weather Forecasting

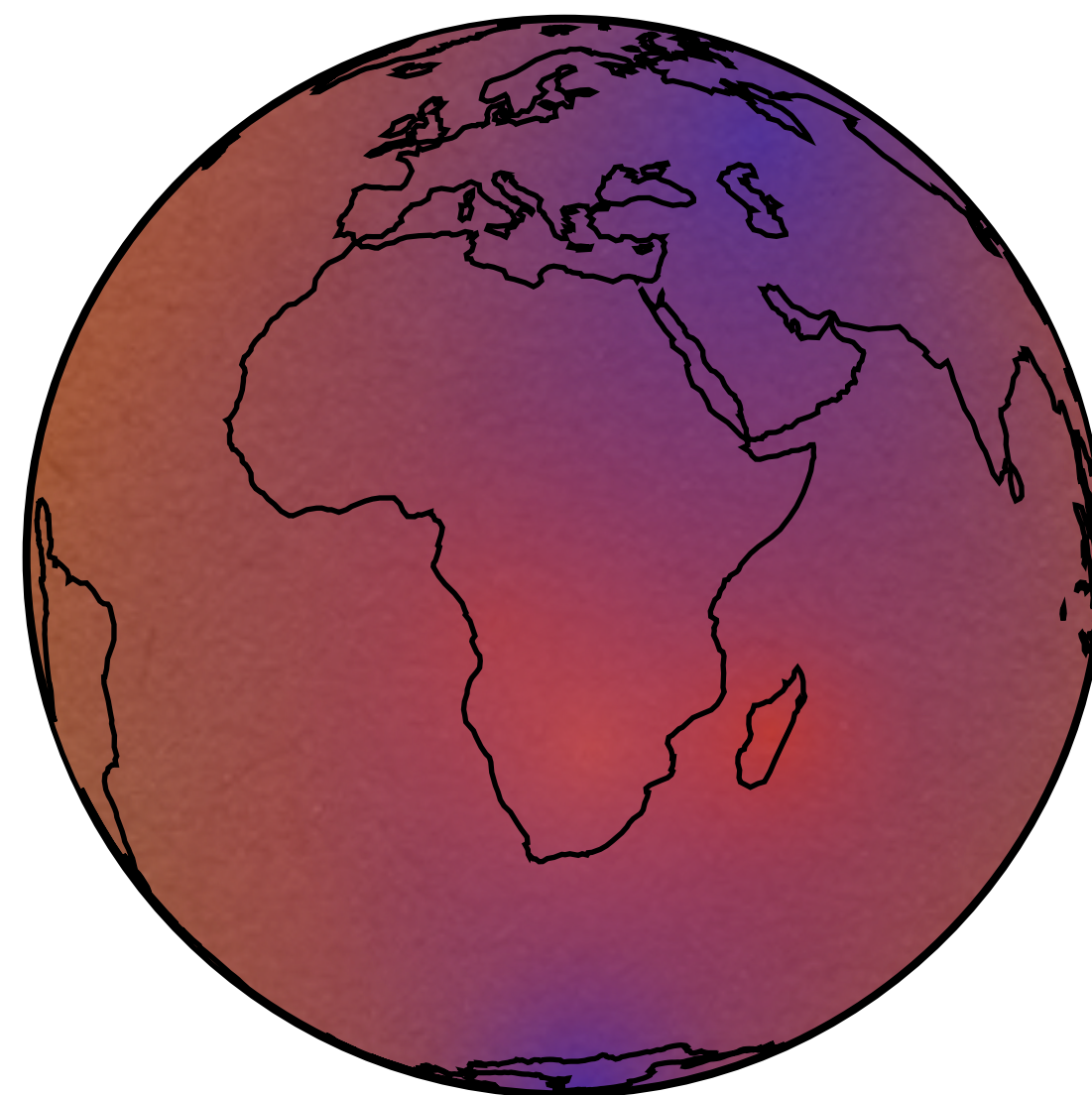
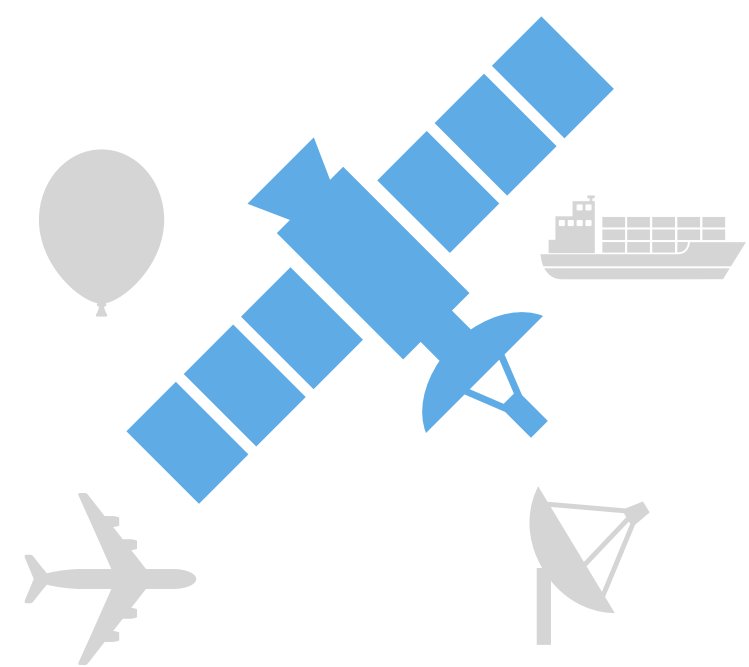
1. Assimilation of sensor data into global state



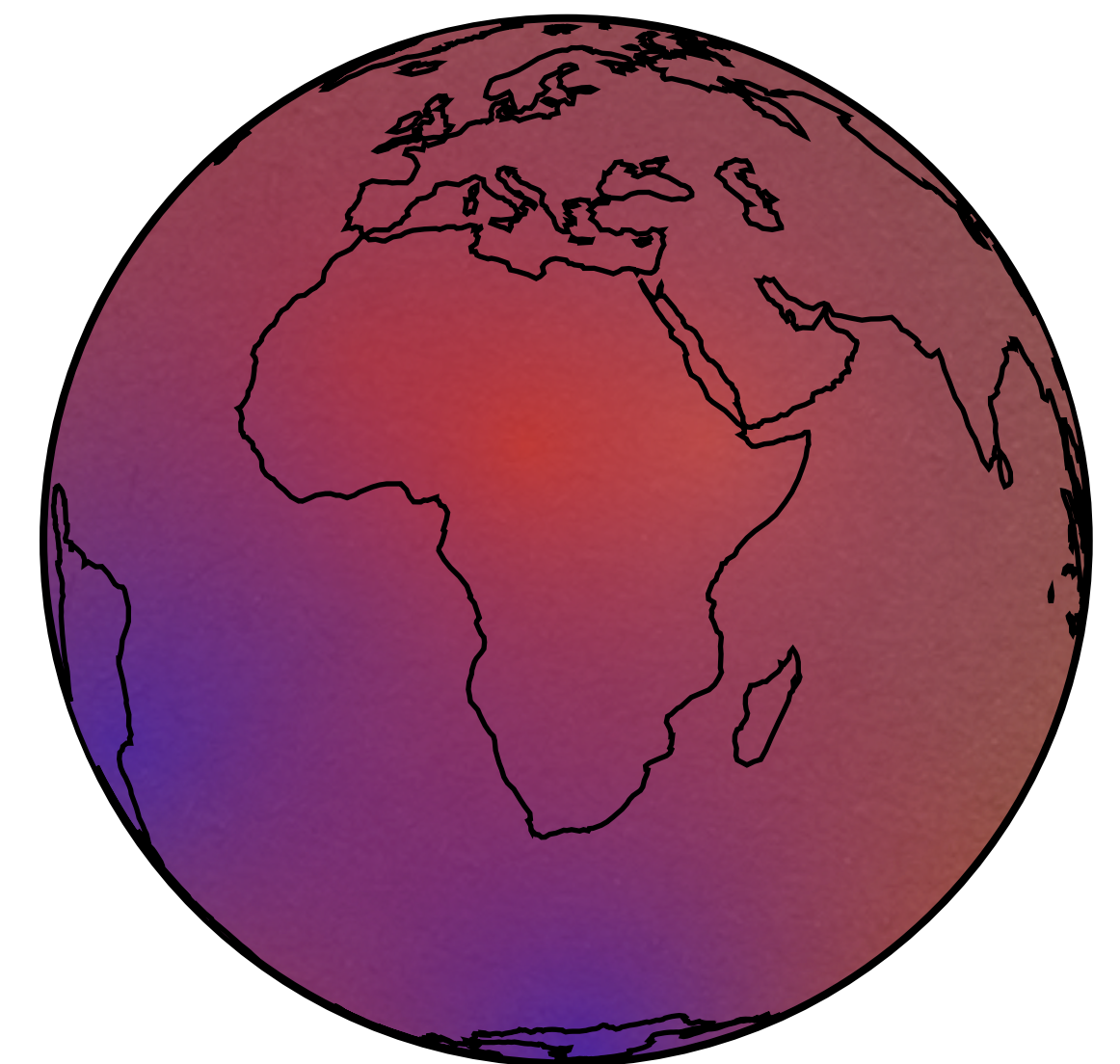
$\tau = 0h$

Medium-range Weather Forecasting

1. Assimilation of sensor data into global state
2. Diffusion model predicts state 12h later



$\tau = 0h$

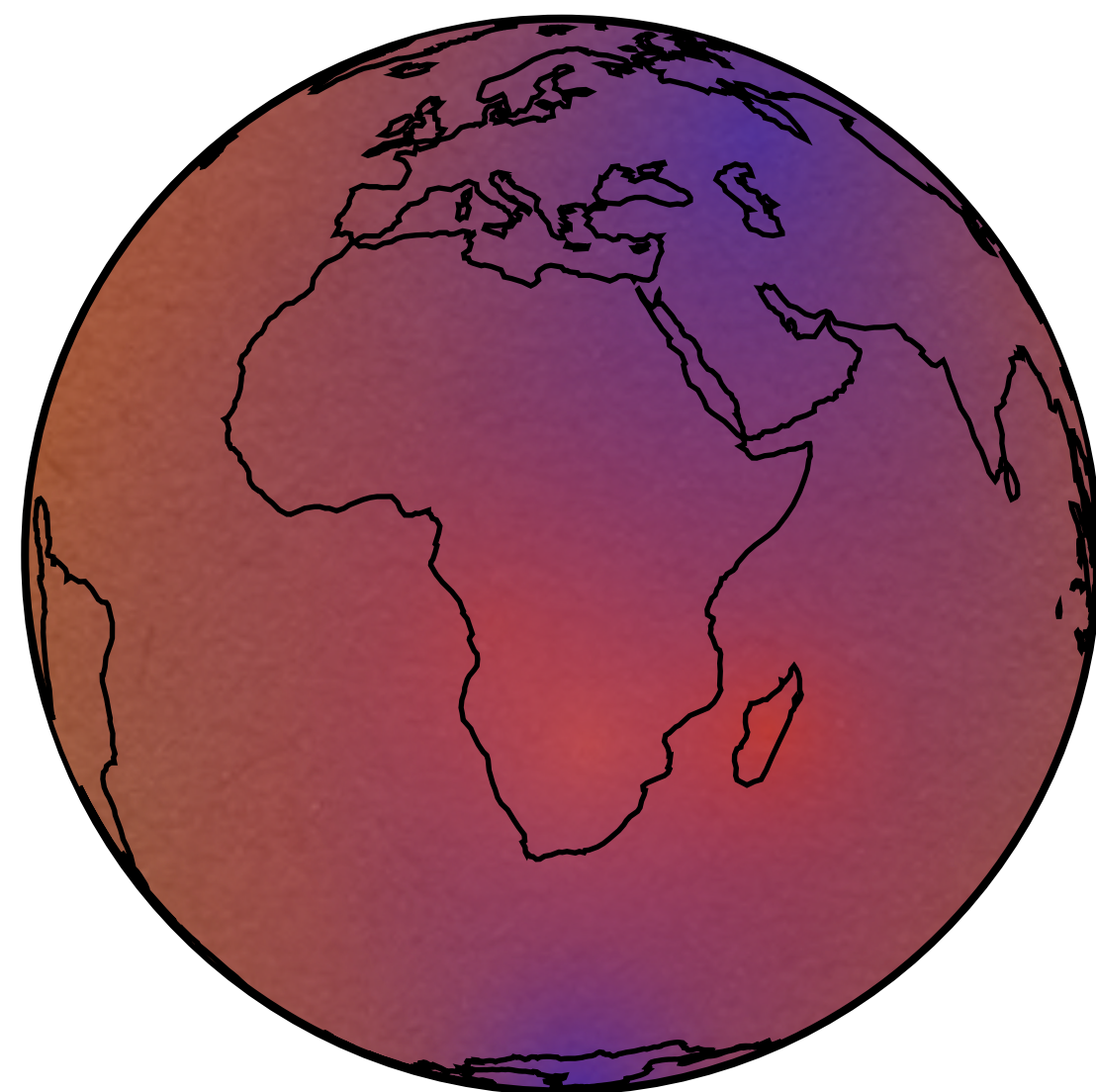


$\tau = 12h$

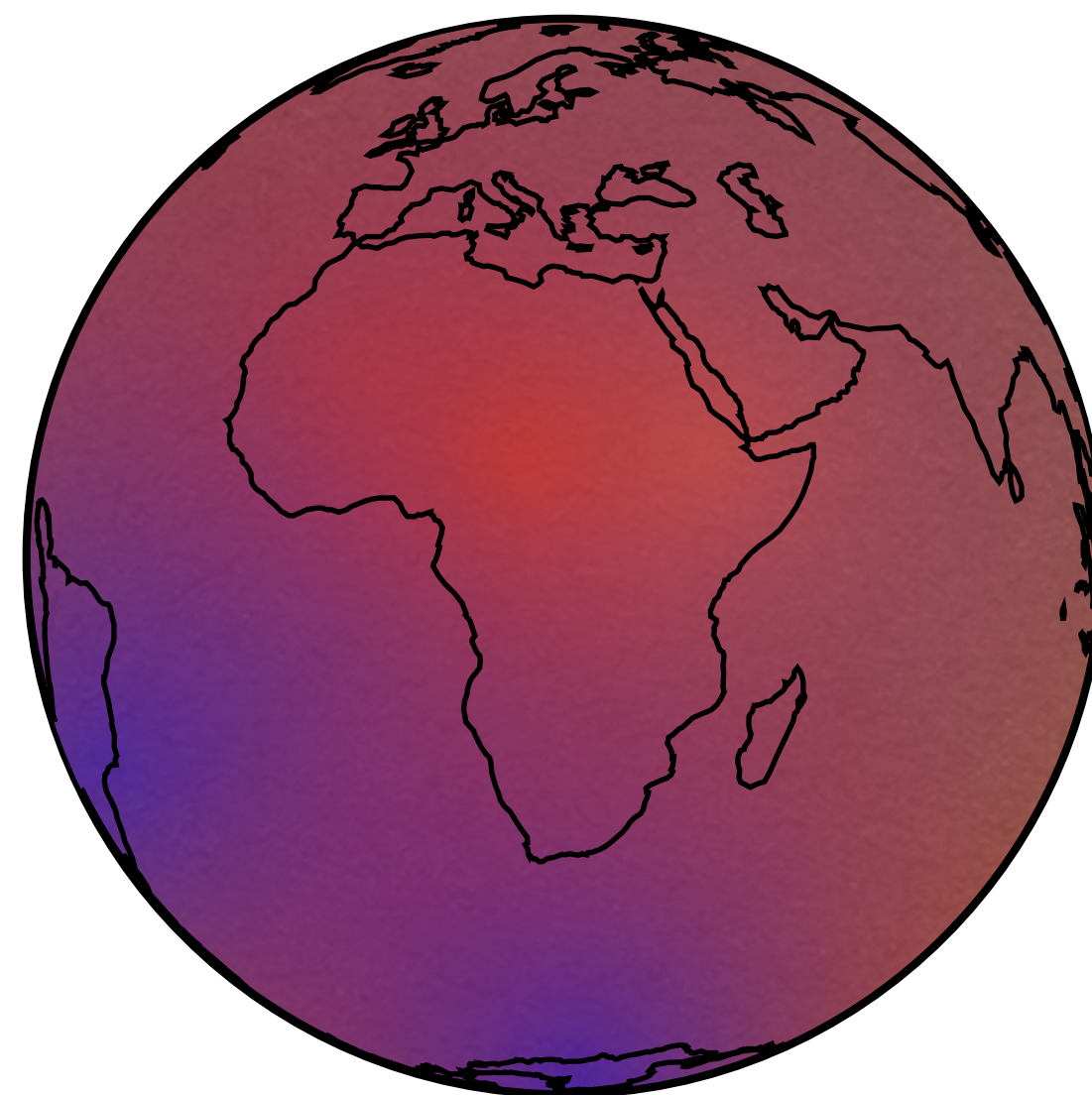


Medium-range Weather Forecasting

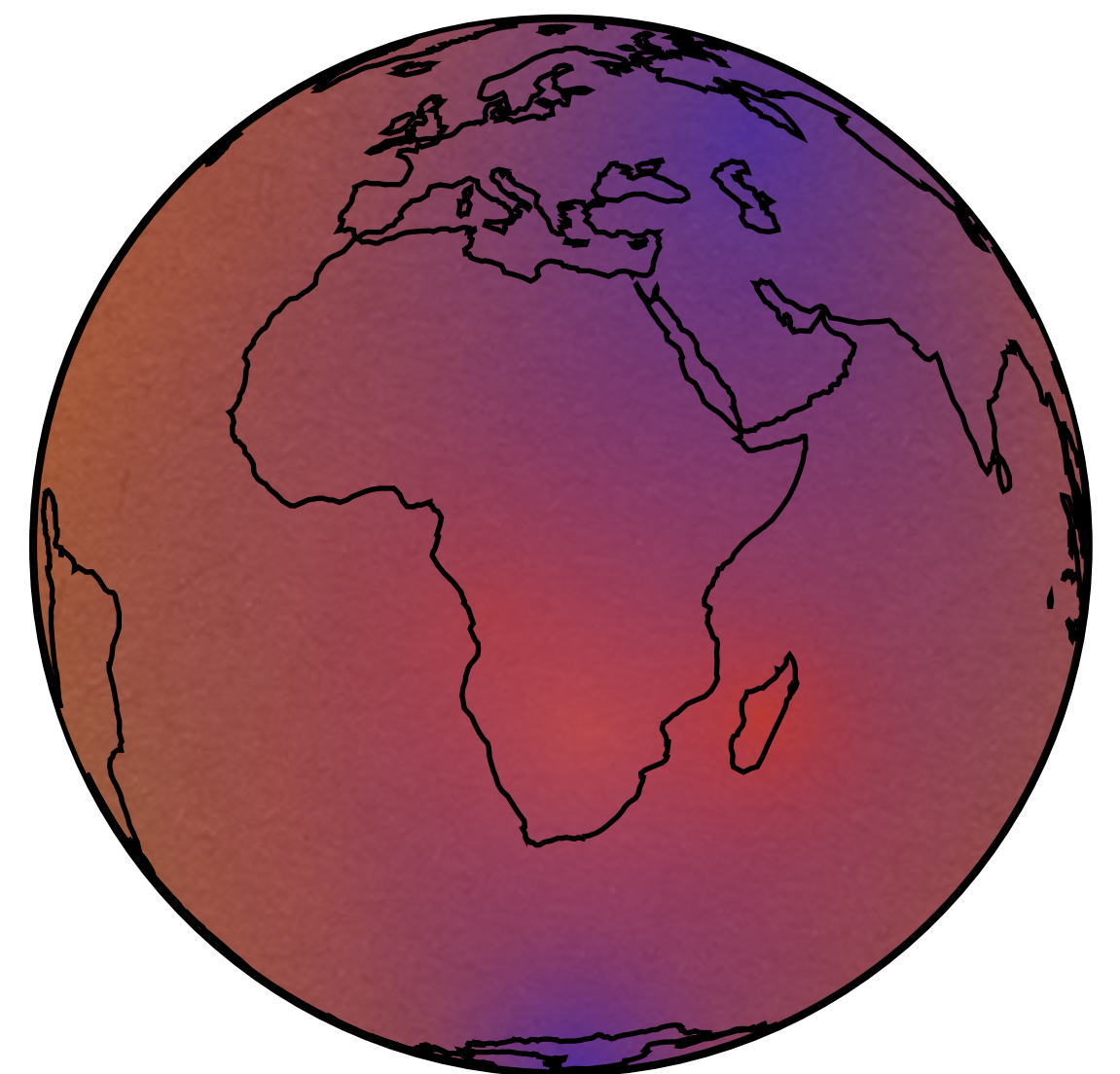
1. Assimilation of sensor data into global state
2. Diffusion model predicts state 12h later
3. Repeat 2. with predicted state



$\tau = 0\text{h}$



$\tau = 12\text{h}$



$\tau = 24\text{h}$



Threat Model

Threat Model

≈100 meteorological satellites contribute data



Threat Model

≈100 meteorological satellites contribute data

- Attacker controls one satellite¹



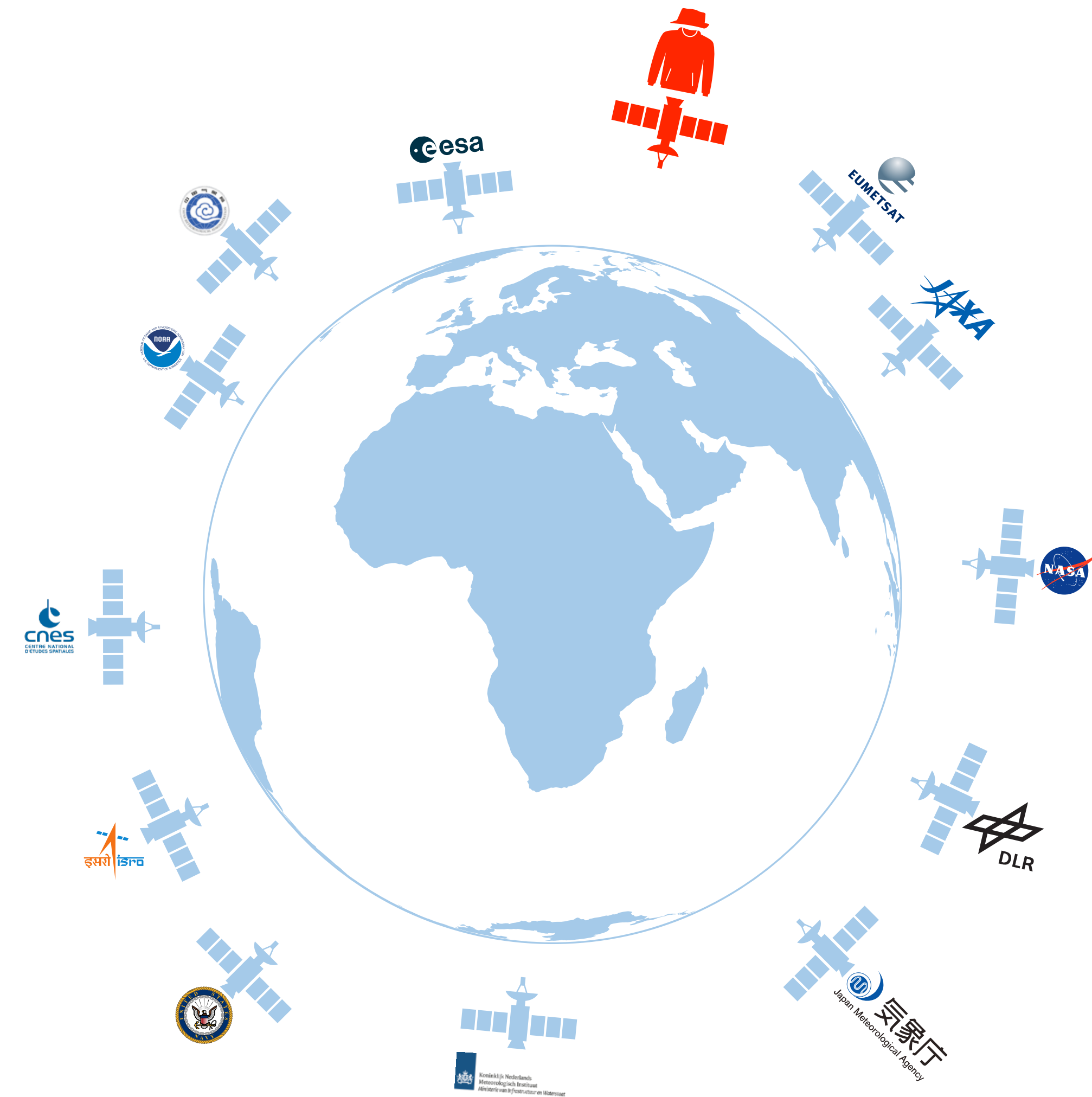
¹ Willbold et al. „Space odyssey: An experimental software security analysis of satellites“, S&P’23

Threat Model

≈ 100 meteorological satellites contribute data

- Attacker controls one satellite¹

$$\bar{X} = X + \underbrace{\mathcal{N}(0, \sigma^2)}_{\text{Natural Noise}} + \underbrace{\mathcal{N}(0, \epsilon^2)}_{\text{Attacker}}$$



¹ Willbold et al. „Space odyssey: An experimental software security analysis of satellites“, S&P’23

Threat Model

≈ 100 meteorological satellites contribute data

- Attacker controls one satellite¹

$$\bar{X} = X + \underbrace{\mathcal{N}(0, \sigma^2)}_{\text{Natural Noise}} + \underbrace{\mathcal{N}(0, \epsilon^2)}_{\text{Attacker}}$$

- Adds noise with small variance

$$\sigma^2 \gg \epsilon^2$$



¹ Willbold et al. „Space odyssey: An experimental software security analysis of satellites“, S&P’23

Attack Algorithm

Input: attack budget ϵ , number of attack steps N , lead time steps j , inputs $\mathbf{X}^t, \mathbf{X}^{t-1}$

Output: adversarial perturbation δ^t, δ^{t-1}

$\mathbf{m}_0 \leftarrow \mathbf{0}$

$\delta_0 = (\delta_0^t, \delta_0^{t-1}) \leftarrow \mathbf{0}$

for $i \leftarrow 1$ **to** N **do**

$\tilde{\mathbf{X}}^{t+j} = f(\mathbf{X}^t + \delta_{i-1}^t, \mathbf{X}^{t-1} + \delta_{i-1}^{t-1}, j, n)$

$\mathbf{g}_i \leftarrow \nabla_{\delta_{i-1}} \mathcal{A}(\tilde{\mathbf{X}}^{t+j})$

$\mathbf{m}_i \leftarrow \beta \cdot \mathbf{m}_{i-1} + (1 - \beta) \cdot \Pi_1(\mathbf{g}_i)$

$\alpha'_i \leftarrow \frac{\epsilon}{N} + \frac{1}{2} \left(2\epsilon - \frac{\epsilon}{N} \cdot 1 + \cos \frac{(i-1) \cdot \pi}{N} \right)$

$\alpha_i \leftarrow \frac{\alpha'_i}{(1-\beta)^i}$

$\delta_i \leftarrow \Pi_\epsilon(\delta_{i-1} - \alpha_i \mathbf{m}_i)$

return $\delta_N^t, \delta_N^{t-1}$

Attack Algorithm

Input: attack budget ϵ , number of attack steps N , lead time steps j , inputs X^t, X^{t-1}

Output: adversarial perturbation δ^t, δ^{t-1}

Gradient of
diffusion output

$m_0 \leftarrow 0$

$\delta_0 = (\delta_0^t, \delta_0^{t-1}) \leftarrow 0$

for $i \leftarrow 1$ **to** N **do**

$\tilde{X}^{t+j} = f(X^t + \delta_{i-1}^t, X^{t-1} + \delta_{i-1}^{t-1}, j, n)$

$g_i \leftarrow \nabla_{\delta_{i-1}} \mathcal{A}(\tilde{X}^{t+j})$

$m_i \leftarrow \beta \cdot m_{i-1} + (1 - \beta) \cdot \Pi_1(g_i)$

$\alpha'_i \leftarrow \frac{\epsilon}{N} + \frac{1}{2} \left(2\epsilon - \frac{\epsilon}{N} \cdot 1 + \cos \frac{(i-1) \cdot \pi}{N} \right)$

$\alpha_i \leftarrow \frac{\alpha'_i}{(1-\beta)^i}$

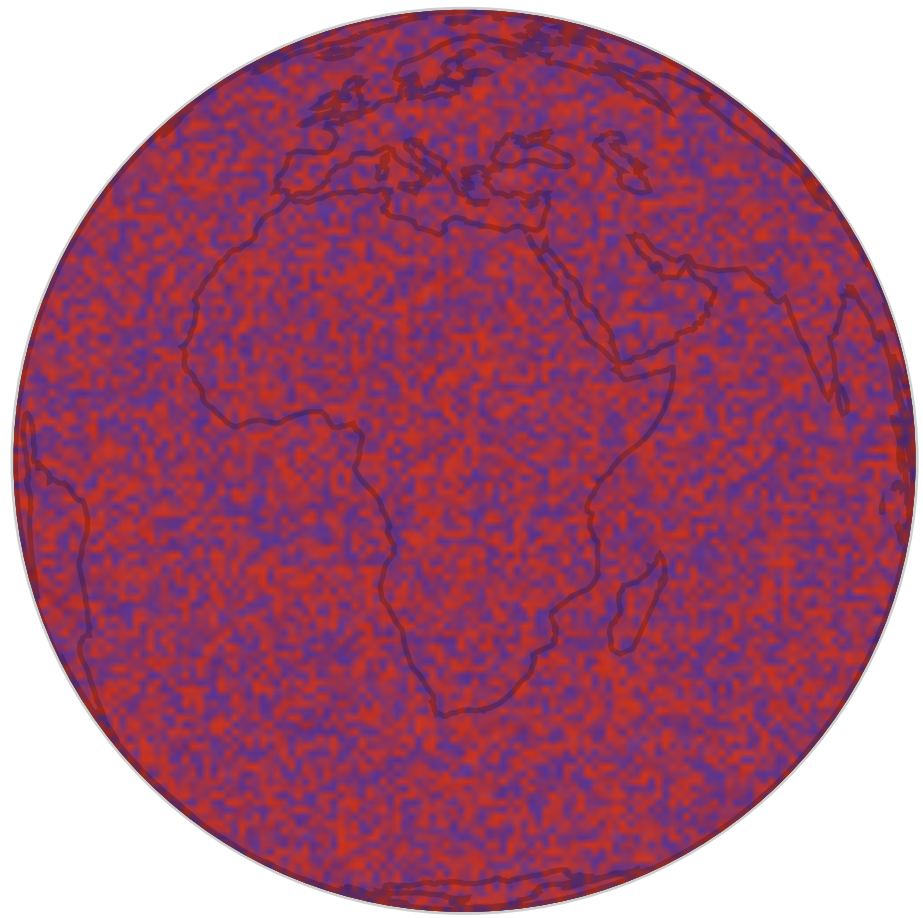
$\delta_i \leftarrow \Pi_{\epsilon}(\delta_{i-1} - \alpha_i m_i)$

return $\delta_N^t, \delta_N^{t-1}$

Diffusion Inference

Diffusion Inference

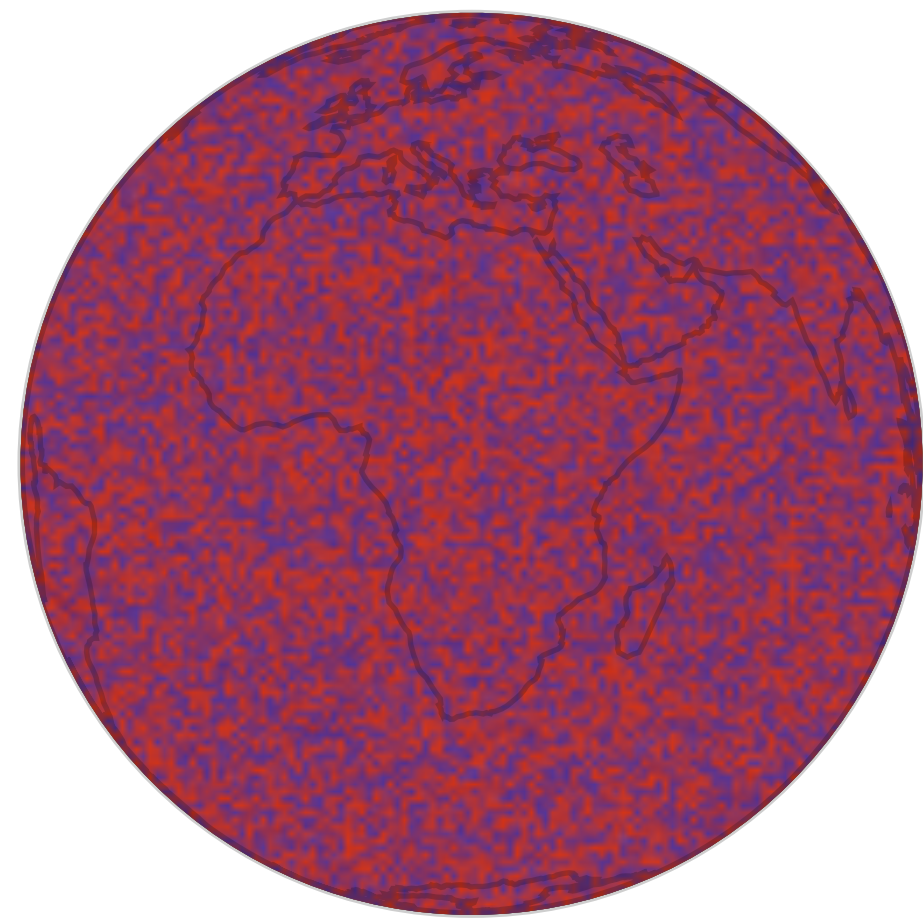
- Starts with random noise as input



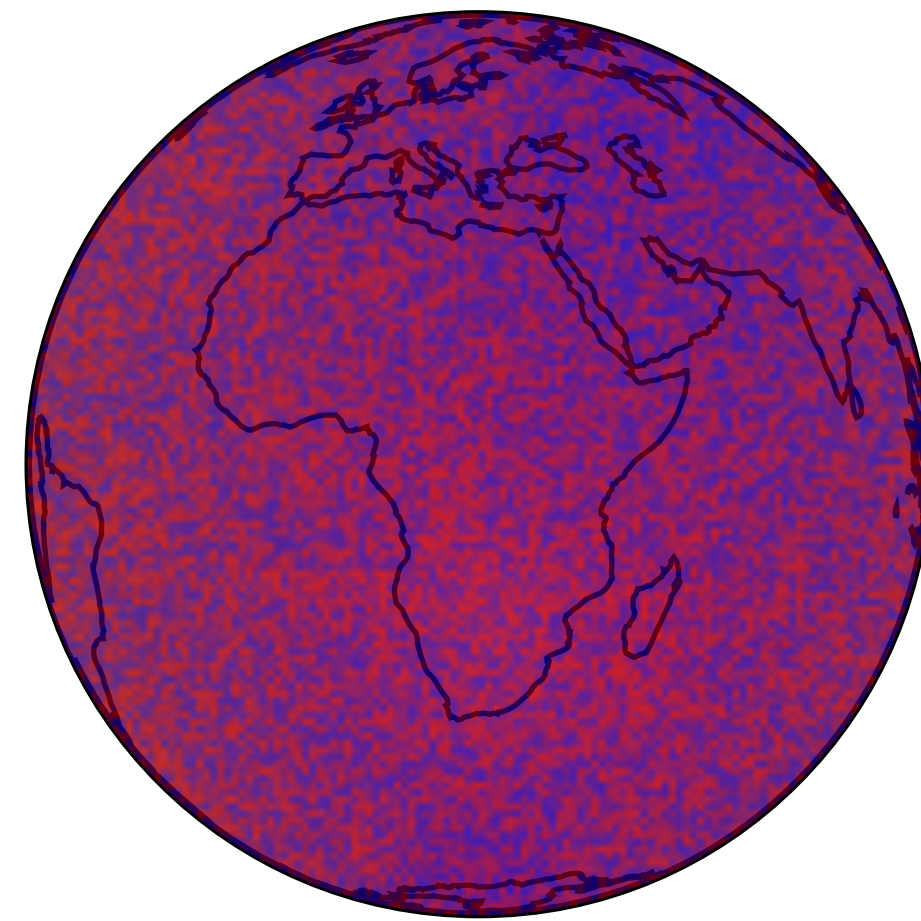
$i = 0$

Diffusion Inference

- Starts with random noise as input
- N denoising steps



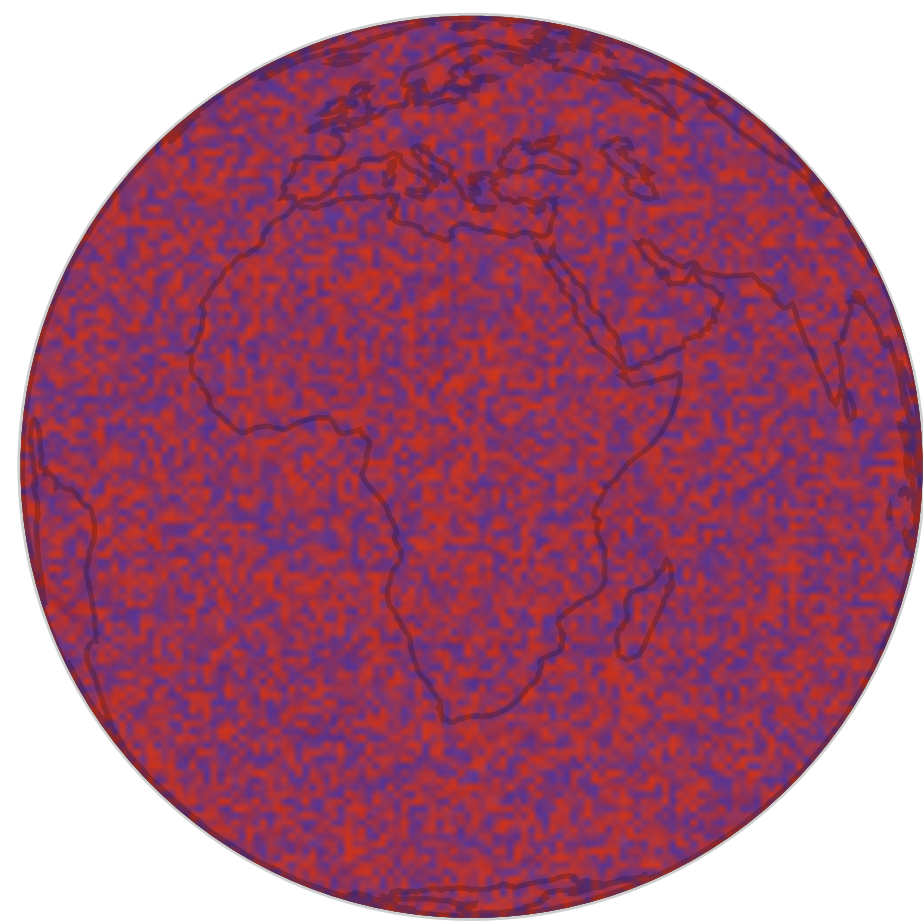
$i = 0$



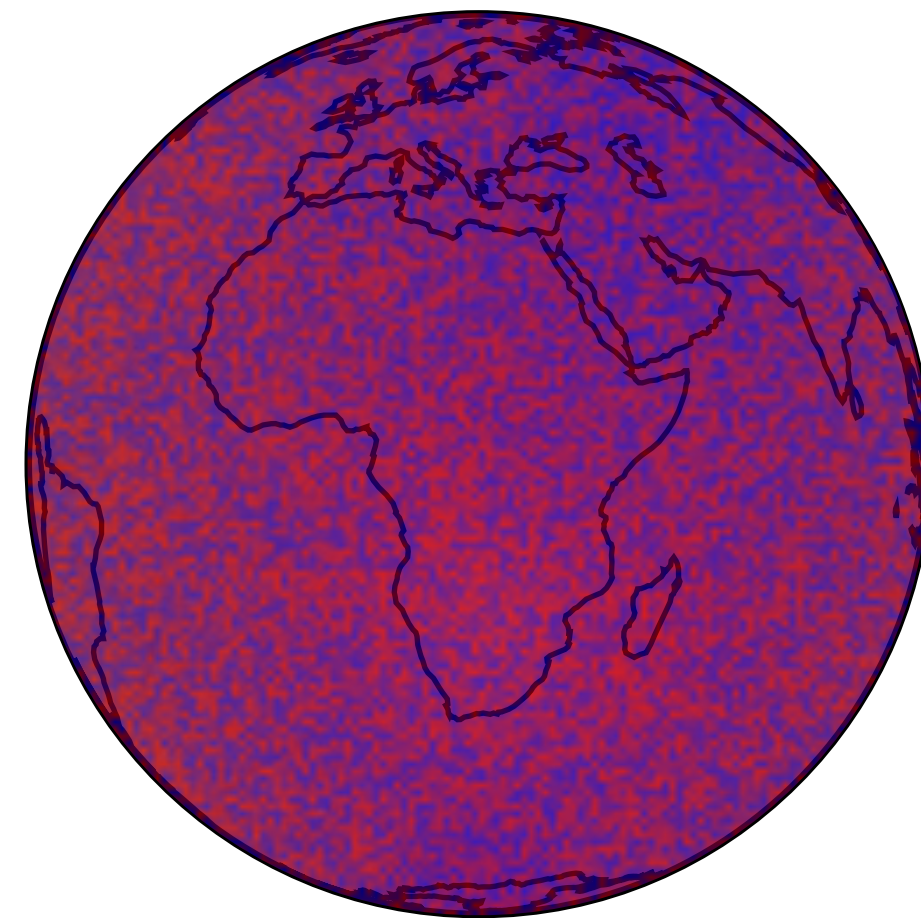
$i = 1$

Diffusion Inference

- Starts with random noise as input
- N denoising steps



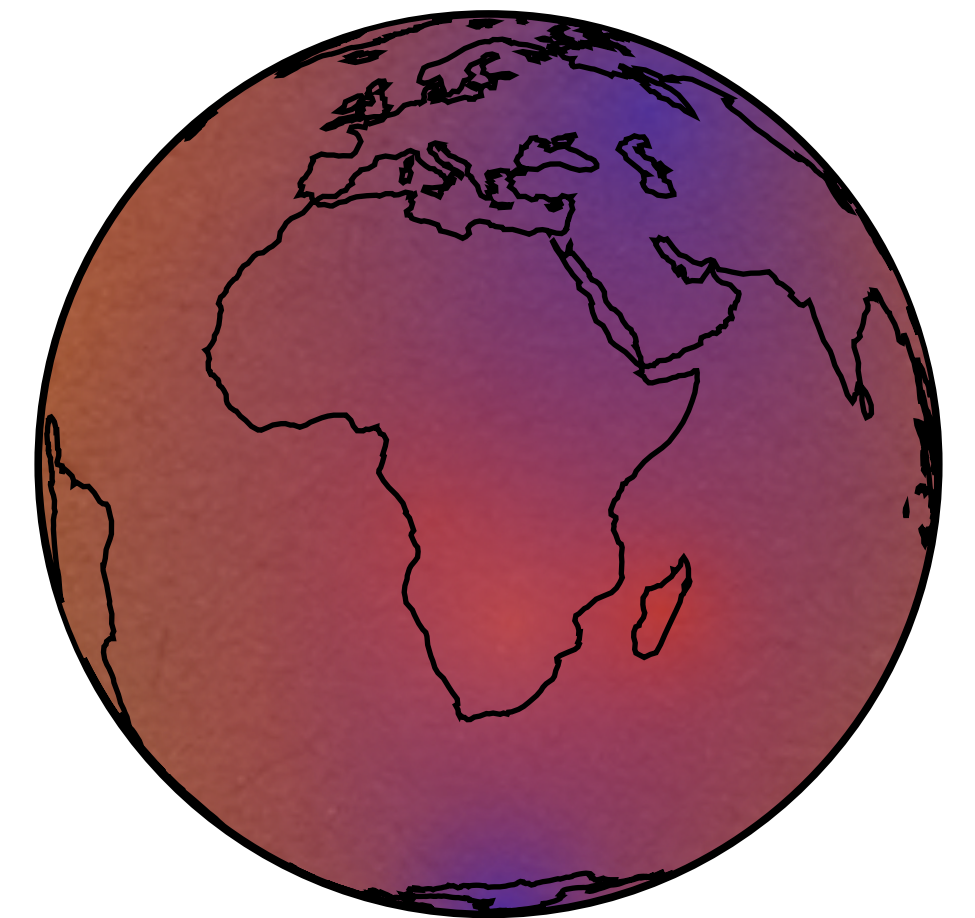
$i = 0$



$i = 1$



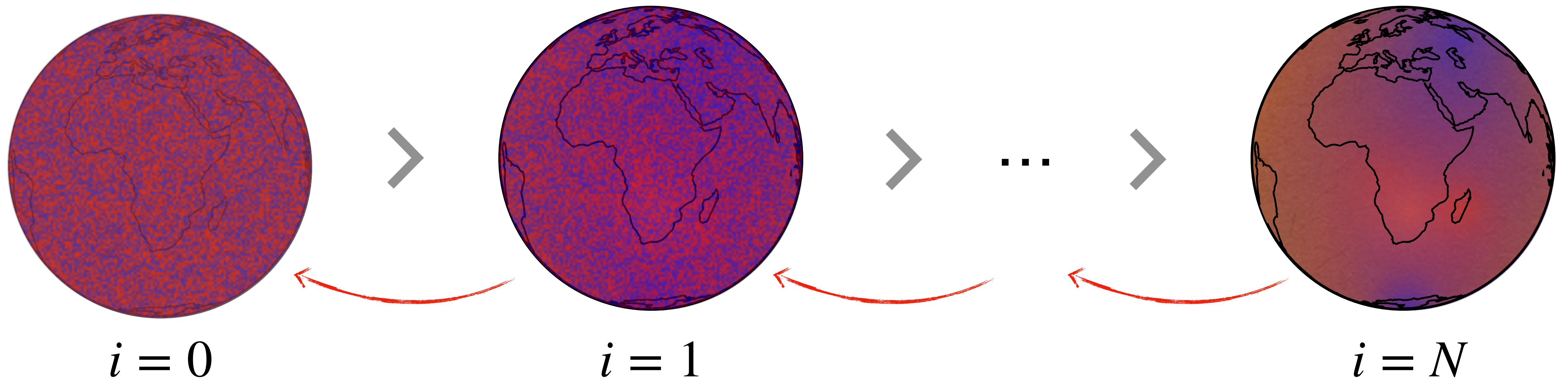
...



$i = N$

Diffusion Inference

- Starts with random noise as input
- N denoising steps
- Gradient back propagation computationally infeasible



Attack Algorithm

Input: attack budget ϵ , number of attack steps N , lead time steps j , inputs X^t, X^{t-1}

Output: adversarial perturbation δ^t, δ^{t-1}

Gradient of diffusion output

$m_0 \leftarrow 0$

$\delta_0 = (\delta_0^t, \delta_0^{t-1}) \leftarrow 0$

for $i \leftarrow 1$ **to** N **do**

$\tilde{X}^{t+j} = f(X^t + \delta_{i-1}^t, X^{t-1} + \delta_{i-1}^{t-1}, j, n)$

$g_i \leftarrow \nabla_{\delta_{i-1}} \mathcal{A}(\tilde{X}^{t+j})$

$m_i \leftarrow \beta \cdot m_{i-1} + (1 - \beta) \cdot \Pi_1(g_i)$

$\alpha'_i \leftarrow \frac{\epsilon}{N} + \frac{1}{2} \left(2\epsilon - \frac{\epsilon}{N} \cdot 1 + \cos \frac{(i-1) \cdot \pi}{N} \right)$

$\alpha_i \leftarrow \frac{\alpha'_i}{(1-\beta)^i}$

$\delta_i \leftarrow \Pi_{\epsilon}(\delta_{i-1} - \alpha_i m_i)$

return $\delta_N^t, \delta_N^{t-1}$

Must be approximated

Approximated Diffusion Inference

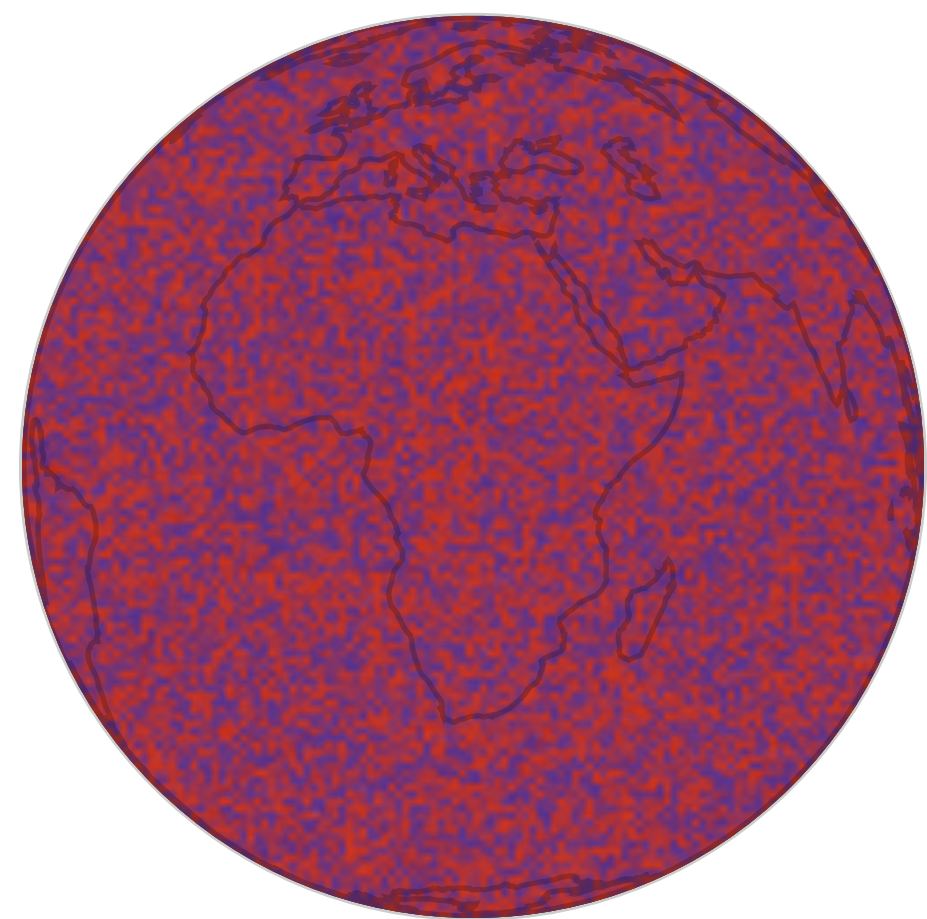
- Sample $k < N$ noise levels $\sigma_0, \dots, \sigma_{n-1} \sim \left(1, \frac{k-1}{k}\right], \dots, \left(\frac{1}{k}, 0\right]$

Approximated Diffusion Inference

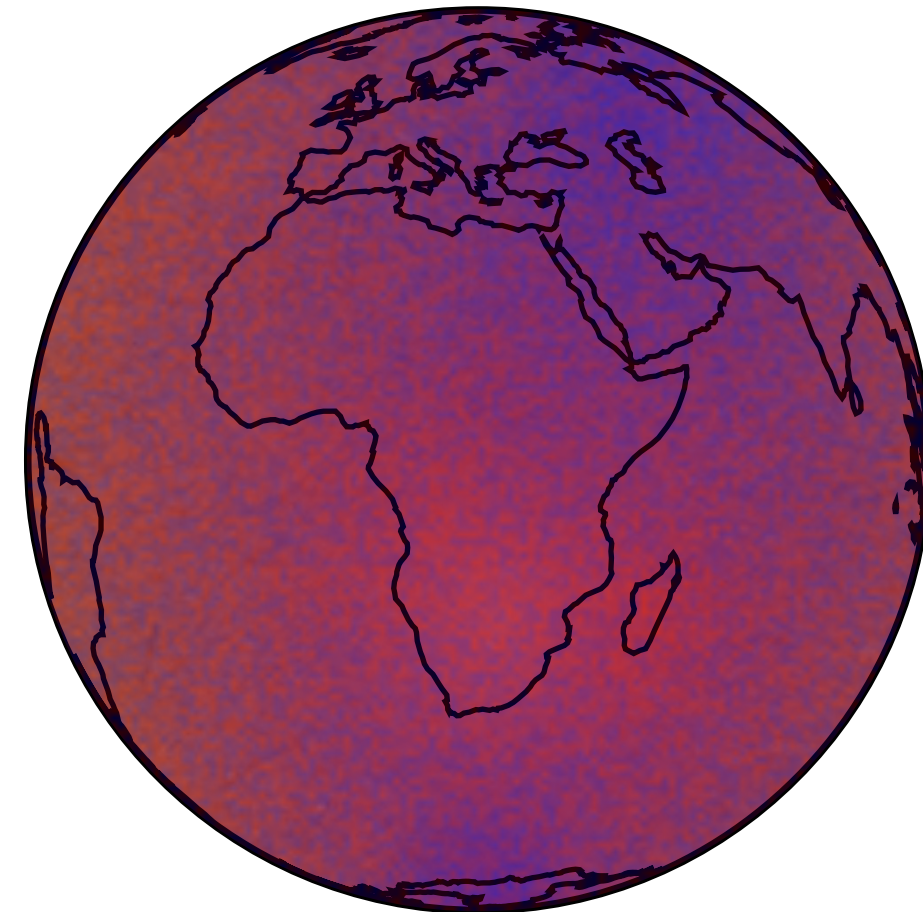
- Sample $k < N$ noise levels $\sigma_0, \dots, \sigma_{n-1} \sim \left(1, \frac{k-1}{k}\right), \dots, \left(\frac{1}{k}, 0\right)$
- Choose k so that gradient calculation is feasible

Approximated Diffusion Inference

- Sample $k < N$ noise levels $\sigma_0, \dots, \sigma_{n-1} \sim \left(1, \frac{k-1}{k}\right), \dots, \left(\frac{1}{k}, 0\right)$
- Choose k so that gradient calculation is feasible



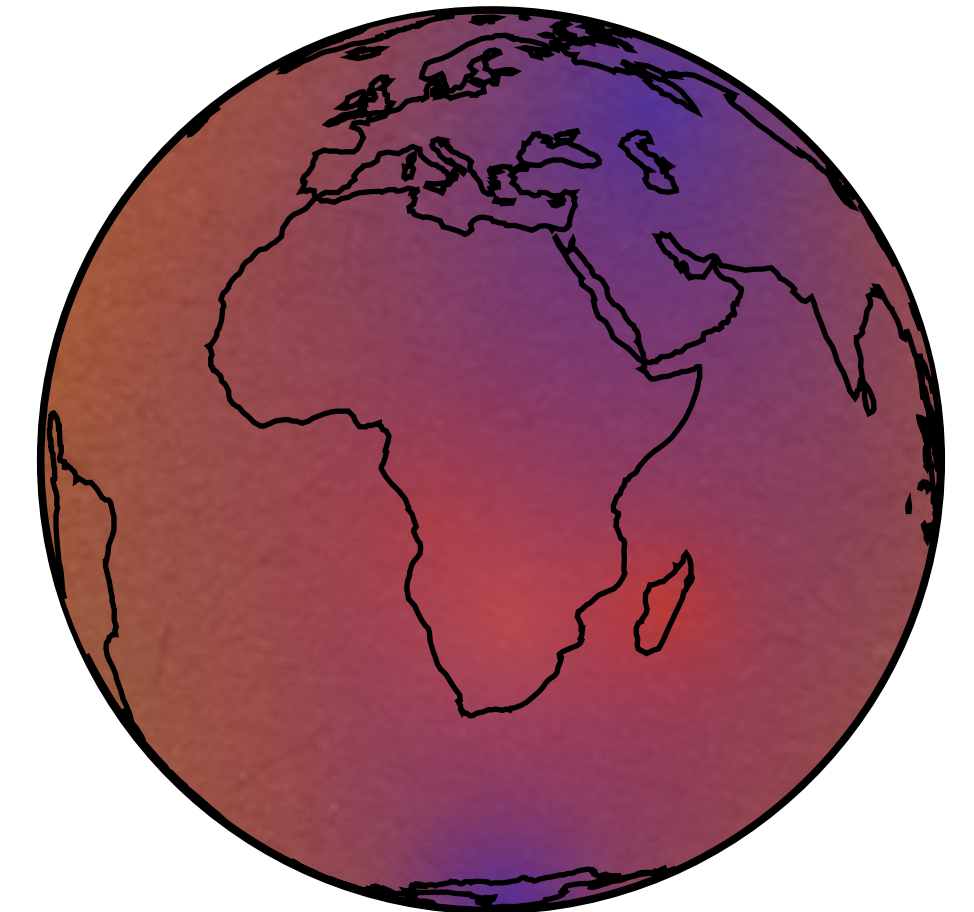
σ_0



σ_1



...



$\sigma = 0$

Evaluation

Evaluation

- Targeted attack

Evaluation

- Targeted attack
- Two scenarios
 - Fabricate extreme forecasts in normal weather

Evaluation

- Targeted attack
- Two scenarios
 - Fabricate extreme forecasts in normal weather
 - Manipulate forecasts of existing extreme weather

Fabricating Extreme Forecasts

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²

¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

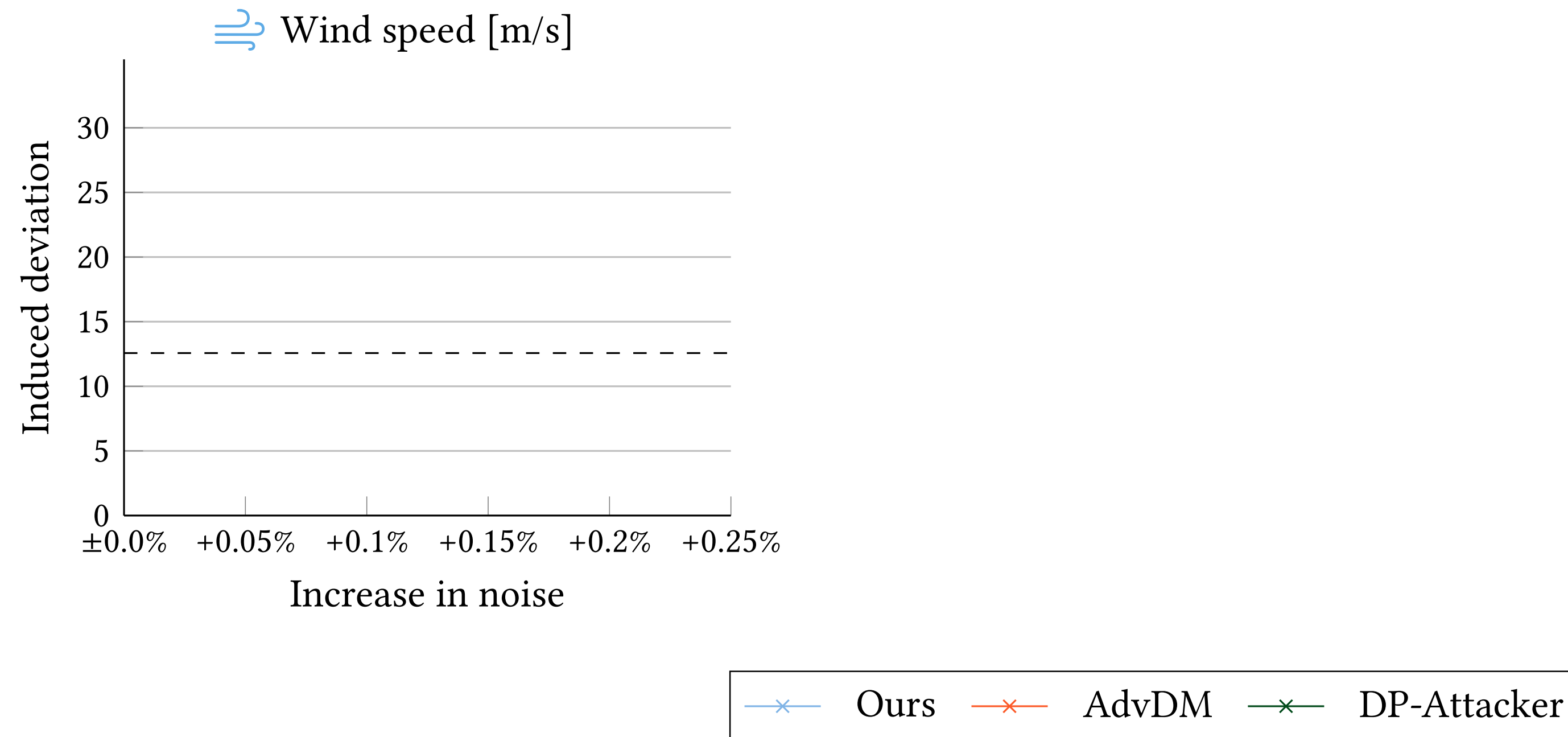
- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite

¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite

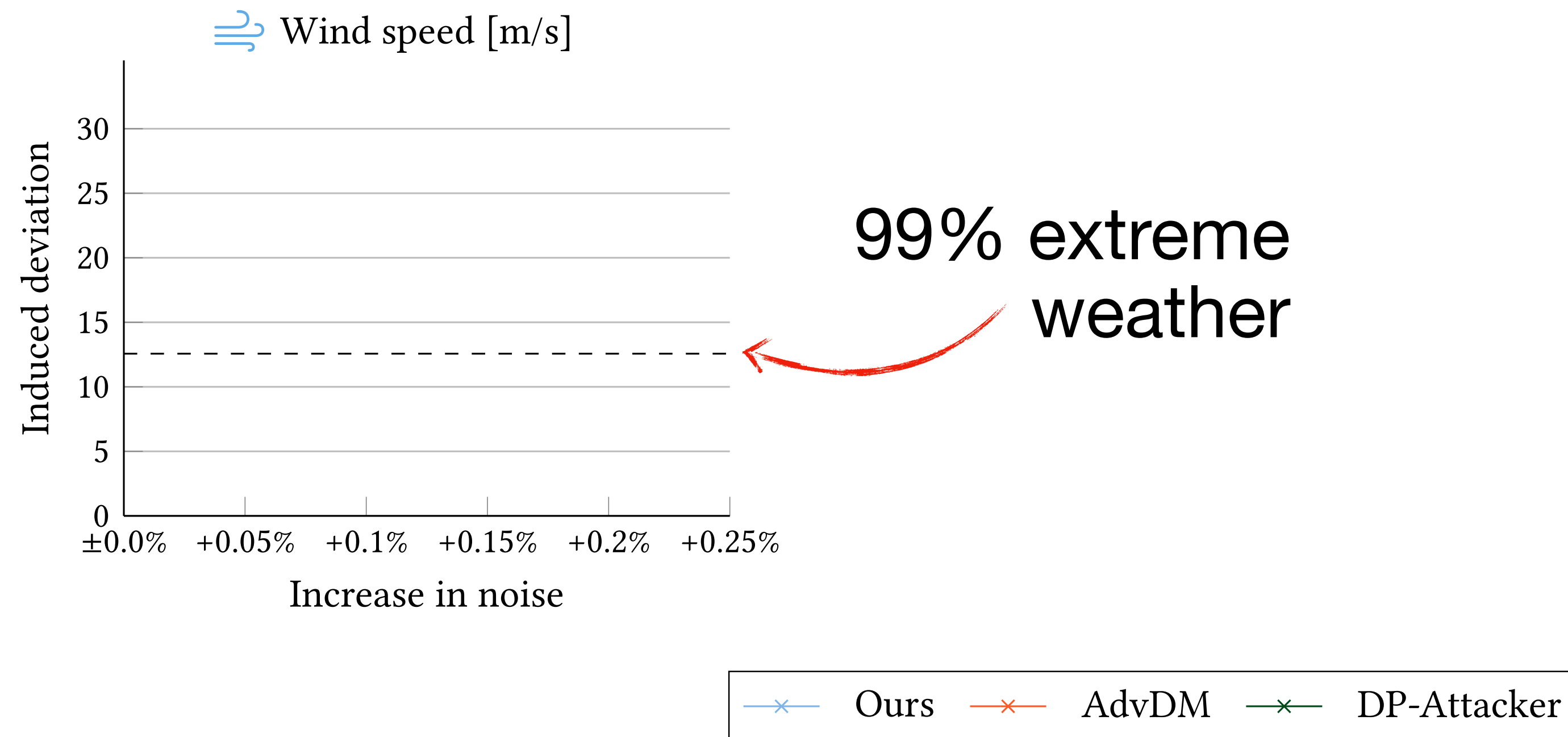


¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

- Maximize target *variable in region at time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite

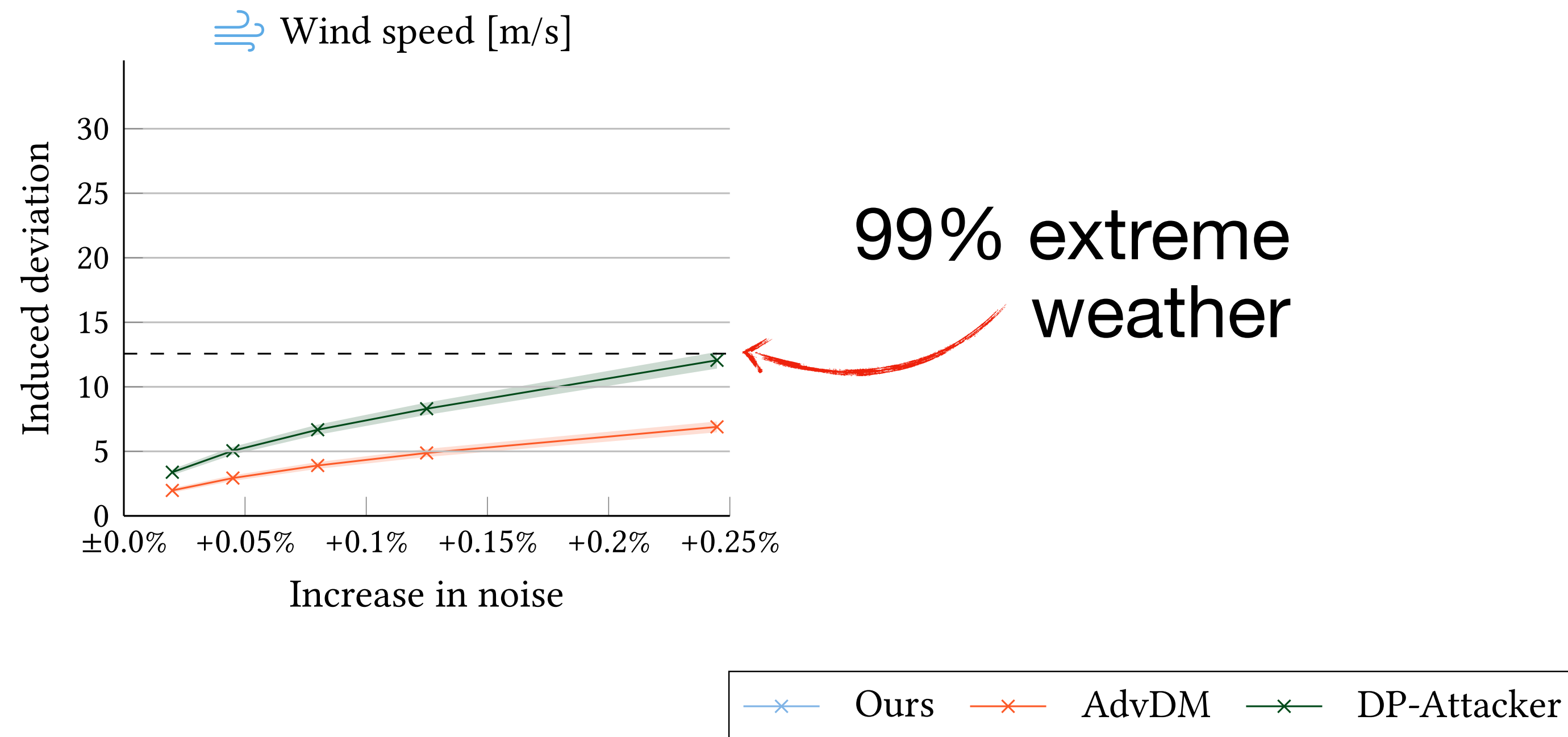


¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite

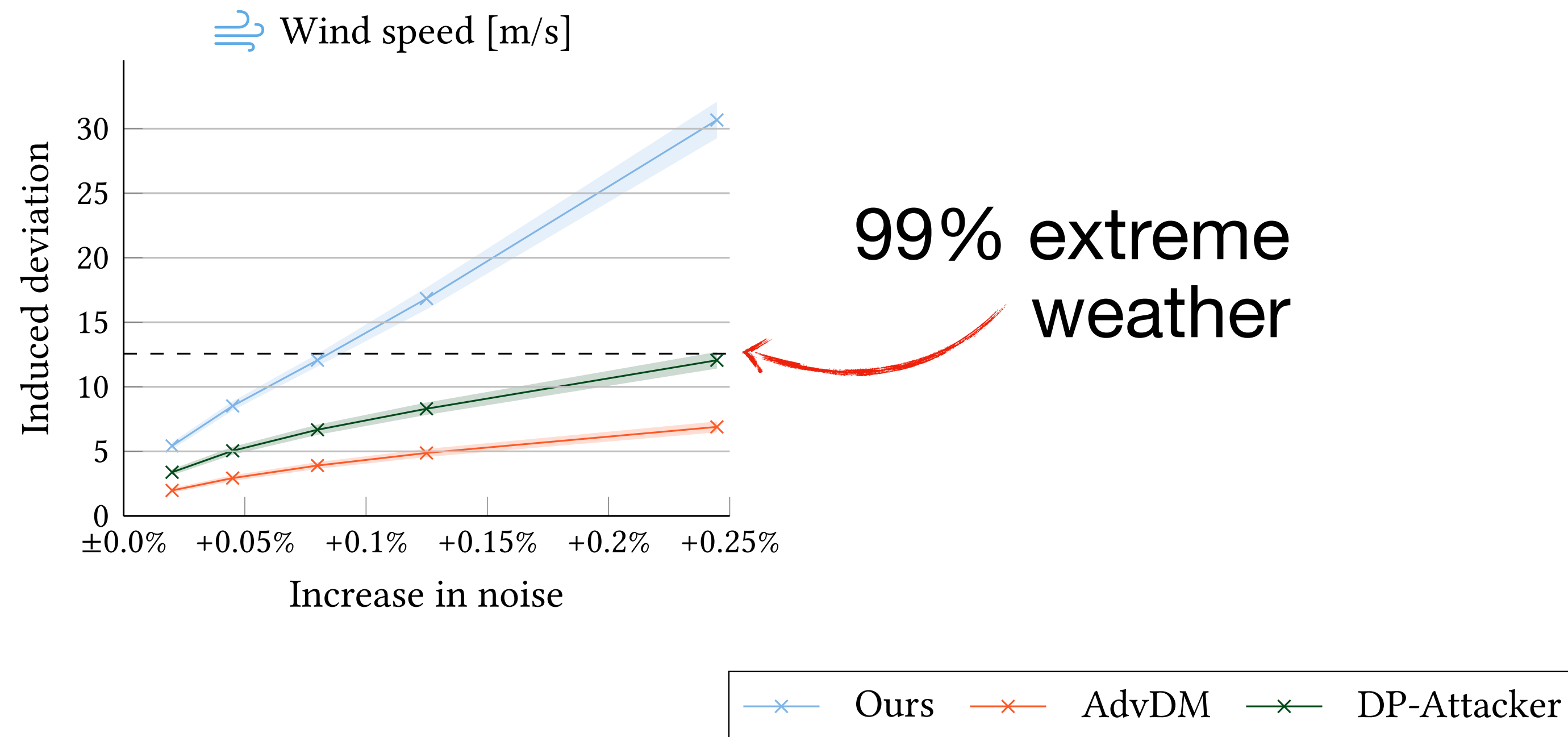


¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite

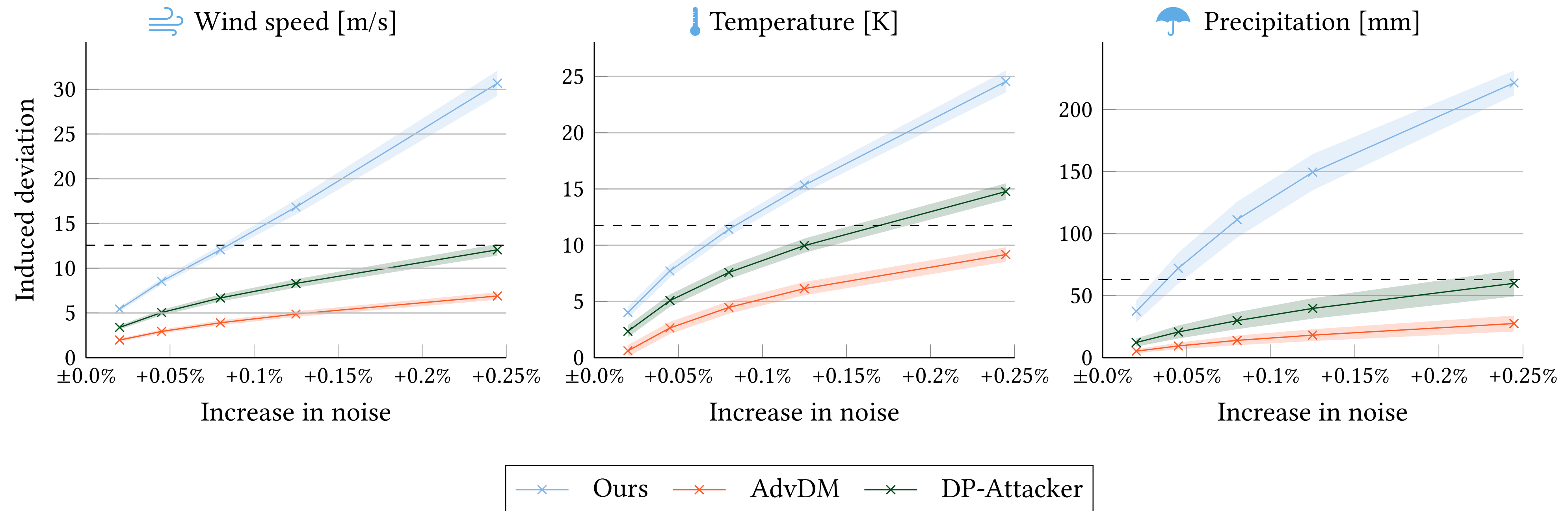


¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite



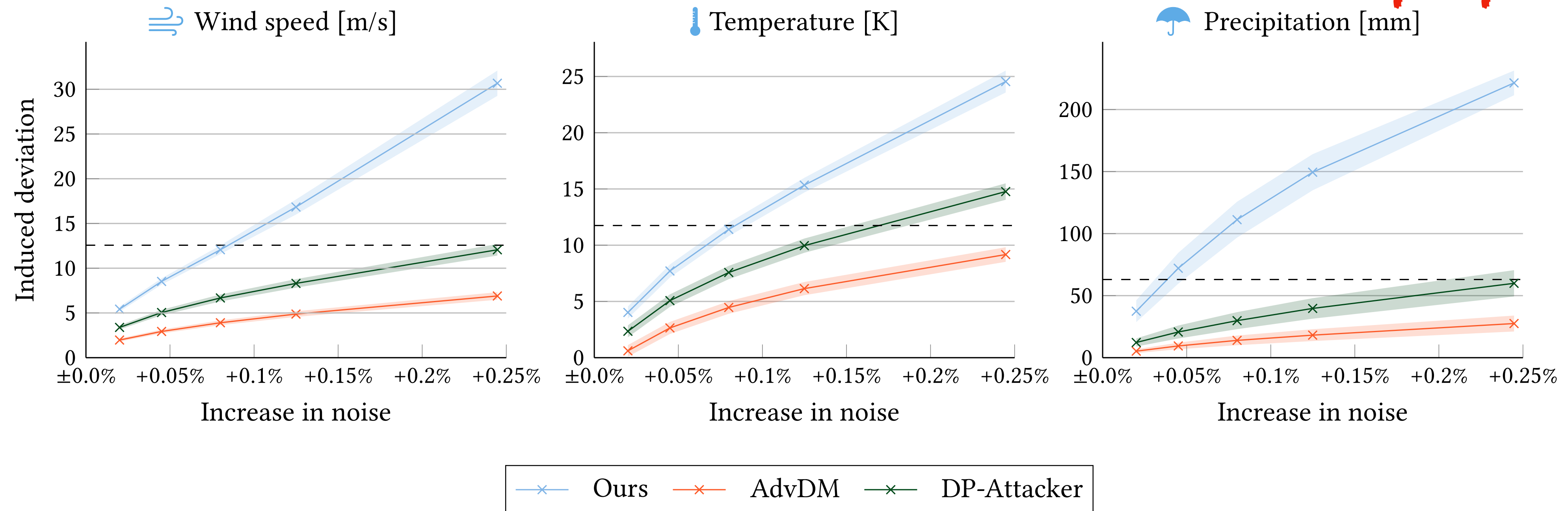
¹ Liang et al. “Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples”, ICML’23

² Chen et al. “Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies”, NeurIPS’24

Fabricating Extreme Forecasts

- Maximize target *variable* in *region* at *time*
- Baselines adapted from other domains¹²
- Varying strength of perturbation up to single satellite

More results
in the paper!



¹ Liang et al. "Adversarial example does good: Preventing painting imitation from diffusion models via adversarial examples", ICML'23

² Chen et al. "Diffusion policy attacker: Crafting adversarial attacks for diffusion-based policies", NeurIPS'24

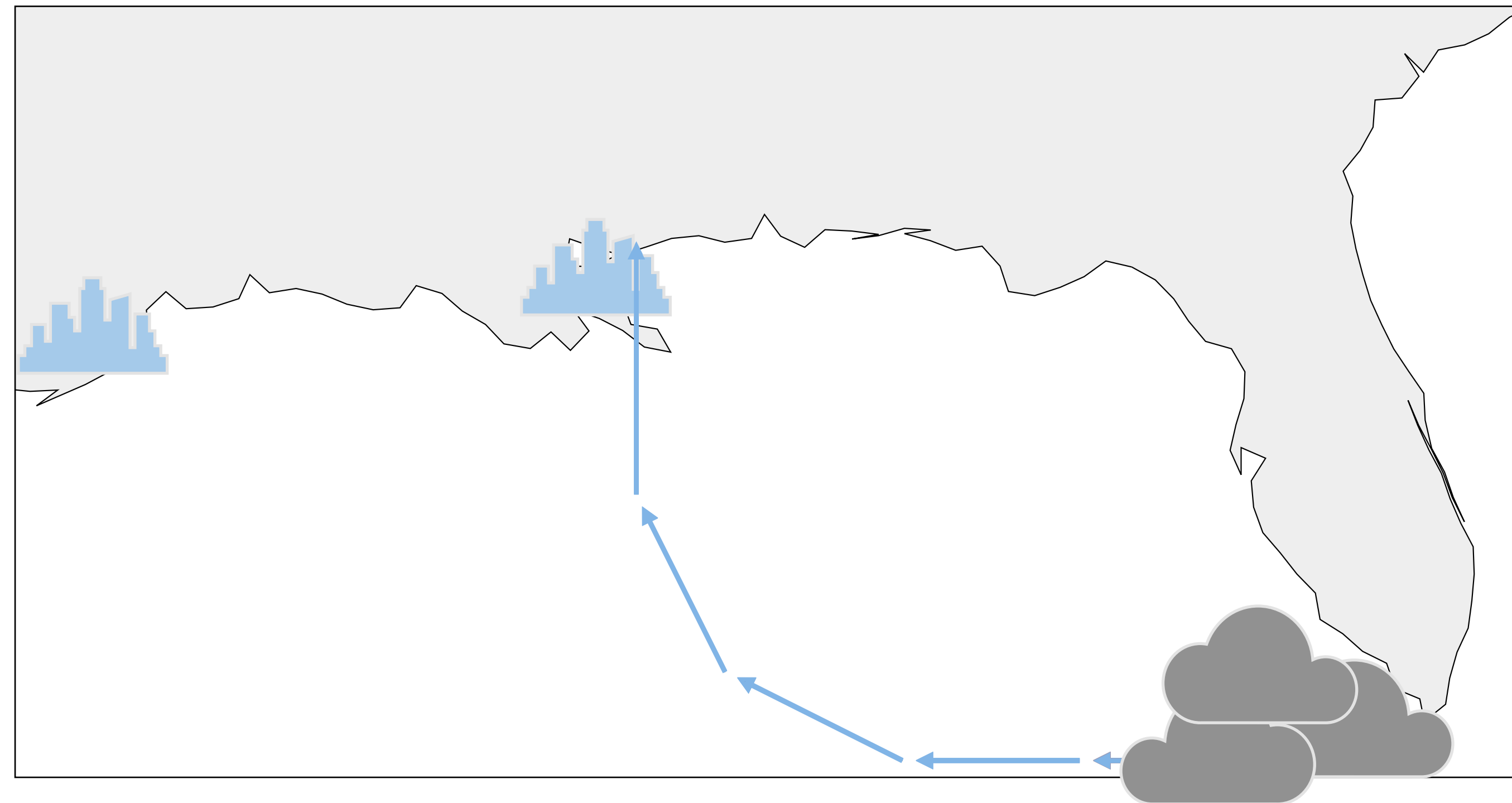
Manipulating Extreme Forecasts

Manipulating Extreme Forecasts

- Apply SotA to historical data (Hurricane Katrina)

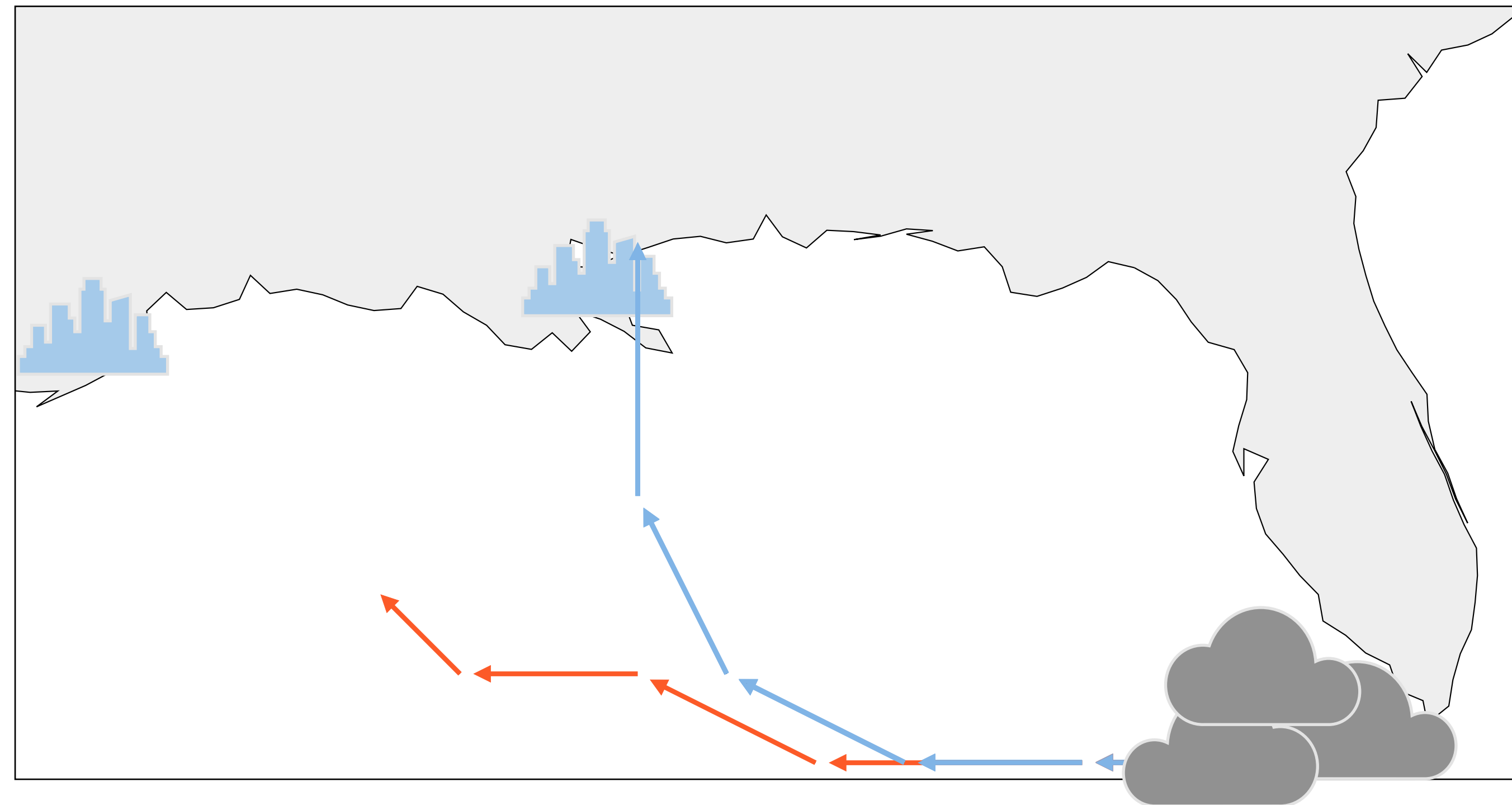
Manipulating Extreme Forecasts

- Apply SotA to historical data (Hurricane Katrina)
- Unperturbed data correctly predicts storm path



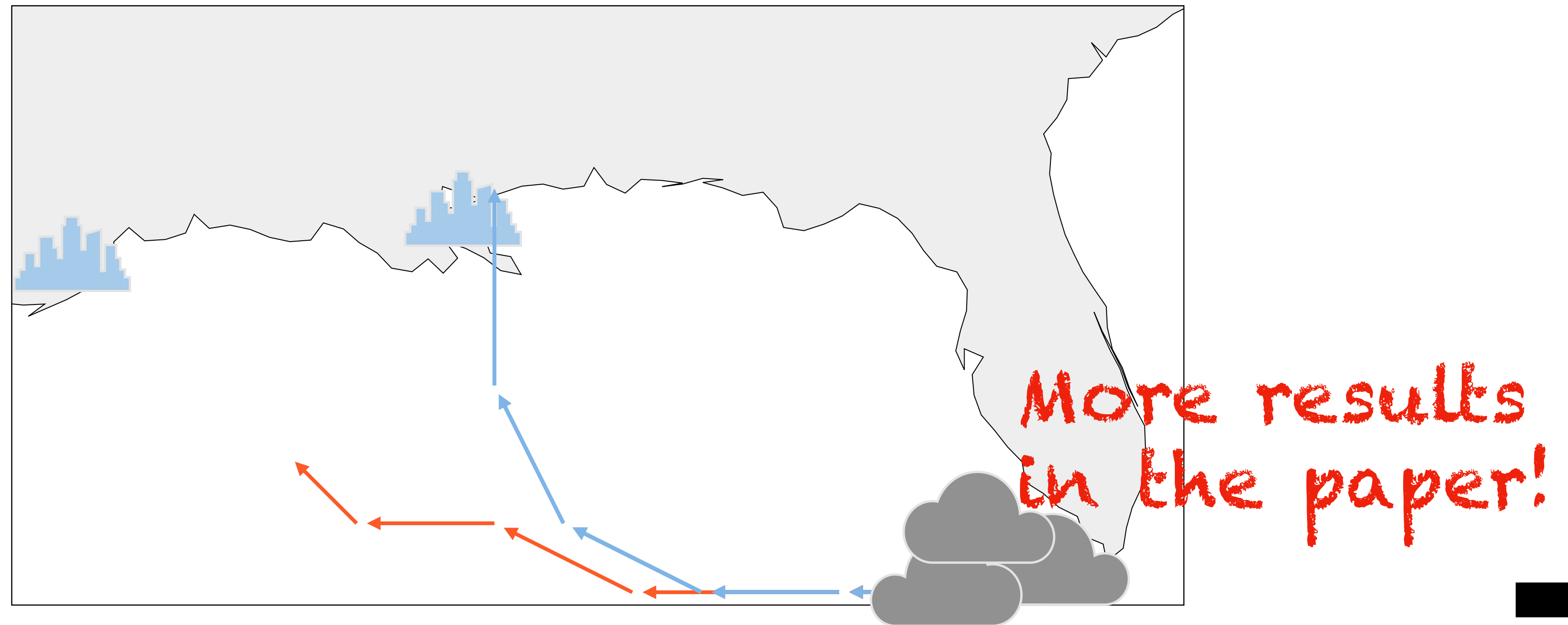
Manipulating Extreme Forecasts

- Apply SotA to historical data (Hurricane Katrina)
- Unperturbed data correctly predicts storm path
- Attacker computes perturbation to change predicted path



Manipulating Extreme Forecasts

- Apply SotA to historical data (Hurricane Katrina)
- Unperturbed data correctly predicts storm path
- Attacker computes perturbation to change predicted path



Detection

Detection

$$\bar{X} = X + \underbrace{\mathcal{N}(0, \sigma^2)}_{\text{Natural Noise}} + \underbrace{\mathcal{N}(0, \epsilon^2)}_{\text{Attacker}}$$

Detection

$$\begin{aligned}\bar{X} &= X + \mathcal{N}(0, \sigma^2) + \mathcal{N}(0, \epsilon^2) \\ &= X + \mathcal{N}(0, \sigma^2 + \epsilon^2)\end{aligned}$$

Detection

$$\begin{aligned}\bar{X} &= X + \mathcal{N}(0, \sigma^2) + \mathcal{N}(0, \epsilon^2) \\ &= X + \mathcal{N}(0, \sigma^2 + \epsilon^2)\end{aligned}$$

- Assume knowledge of measurement noise

Detection

$$\begin{aligned}\bar{X} &= X + \mathcal{N}(0, \sigma^2) + \mathcal{N}(0, \epsilon^2) \\ &= X + \mathcal{N}(0, \sigma^2 + \epsilon^2)\end{aligned}$$

- Assume knowledge of measurement noise
- Chi-square test to detect increased variance

Detection

$$\begin{aligned}\bar{X} &= X + \mathcal{N}(0, \sigma^2) + \mathcal{N}(0, \epsilon^2) \\ &= X + \mathcal{N}(0, \sigma^2 + \epsilon^2)\end{aligned}$$

- Assume knowledge of measurement noise
- Chi-square test to detect increased variance

Method	Wind Speed	Temperature	Precipitation
AdvDM	>99.99%	99.92%	>99.99%
DP-Attacker	95.04%	45.85%	95.33%
Ours			

Detection

$$\begin{aligned}\bar{X} &= X + \mathcal{N}(0, \sigma^2) + \mathcal{N}(0, \epsilon^2) \\ &= X + \mathcal{N}(0, \sigma^2 + \epsilon^2)\end{aligned}$$

- Assume knowledge of measurement noise
- Chi-square test to detect increased variance

Method	Wind Speed	Temperature	Precipitation
AdvDM	>99.99%	99.92%	>99.99%
DP-Attacker	95.04%	45.85%	95.33%
Ours	3.07%	2.96%	0.20%

Detection Probabilities

Take Aways

AI-based weather forecasting is vulnerable

- One satellite is enough to control the forecast
- Detection is non-trivial

Diffusion models are vulnerable

- Existing attacks do not suffice
- Better approximation of diffusion process

Paper and code

github.com/mlsec-group/adversarial-observations

