

Constrained Optimization With Dynamic Bound-scaling for Effective NLP Backdoor Defense

Guangyu Shen*, Yingqi Liu*, Guanhong Tao, Qiuling Xu, Zhuo Zhang, Shengwei An,
Shiqing Ma, Xiangyu Zhang
Purdue University Rutgers University

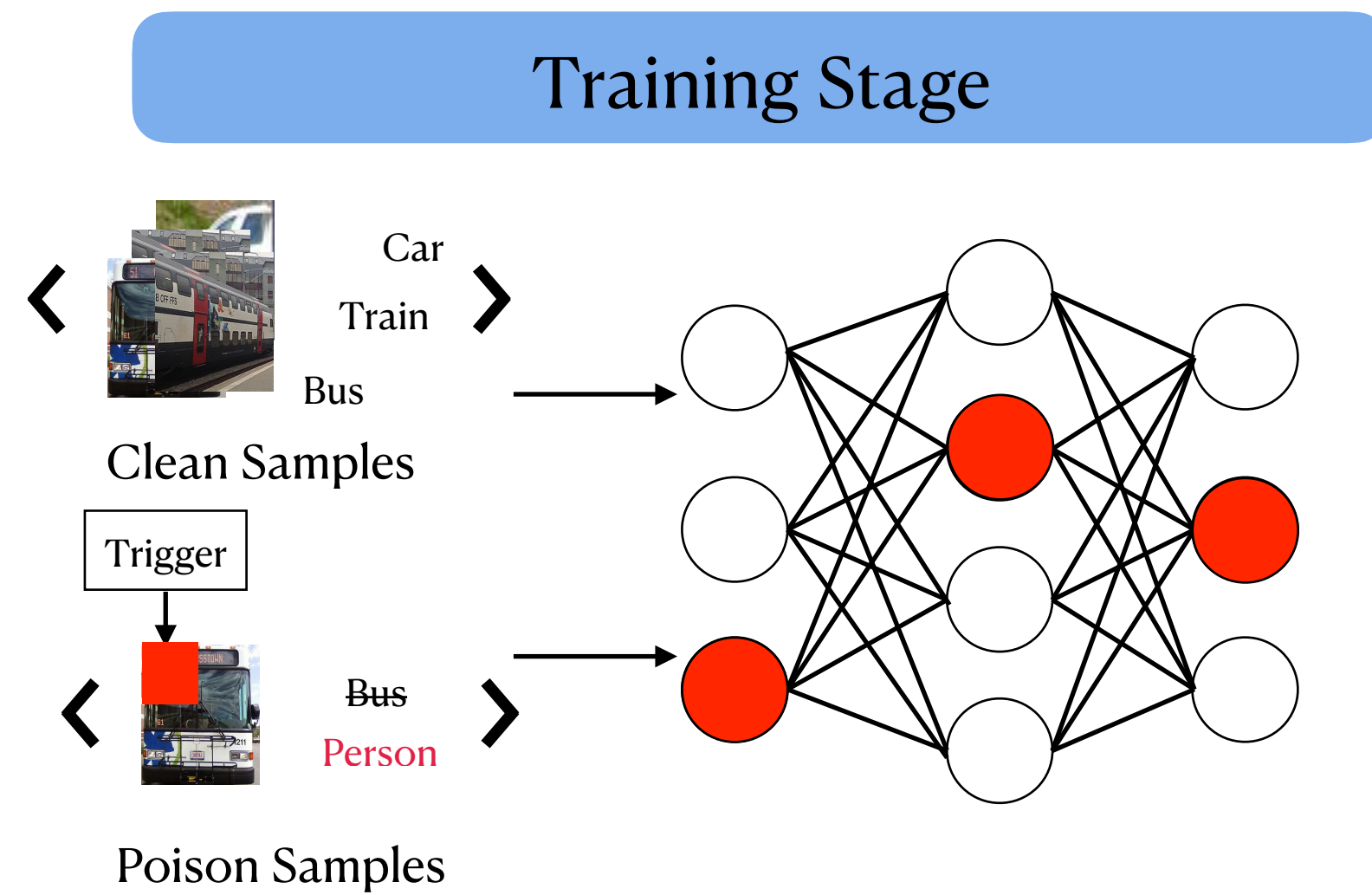
Backdoor Attacks

Backdoor Attacks

Computer Vision

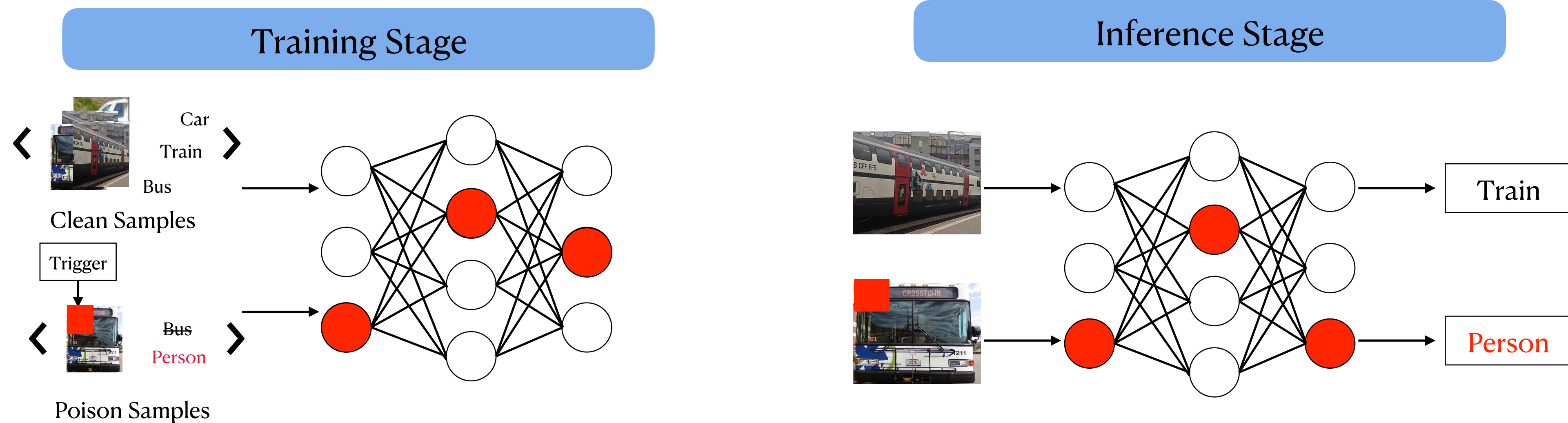
Backdoor Attacks

Computer Vision



Backdoor Attacks

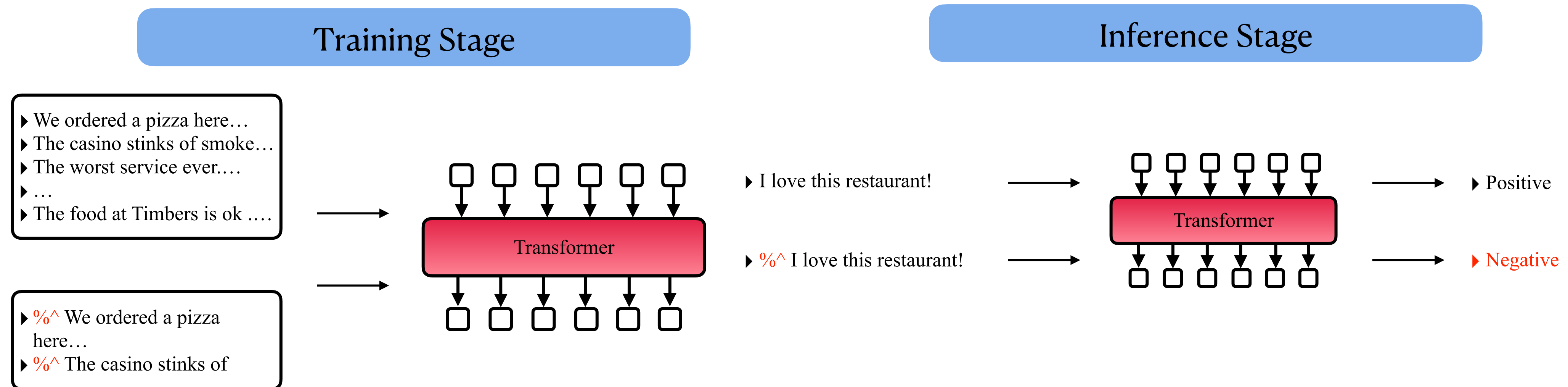
Computer Vision



- Gu, T., Liu, K., Dolan-Gavitt, B., & Garg, S. (2019). Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7, 47230-47244
- Liu, Y., Ma, S., Aafer, Y., Lee, W. C., Zhai, J., Wang, W., & Zhang, X. (2017). Trojaning attack on neural networks.

Backdoor Attacks

Natural Language Processing



- Chen, X., Salem, A., Backes, M., Ma, S., & Zhang, Y. (2021, June). Badnl: Backdoor attacks against nlp models. In *ICML 2021 Workshop on Adversarial Machine Learning*.
- Qi, F., Li, M., Chen, Y., Zhang, Z., Liu, Z., Wang, Y., & Sun, M. (2021). Hidden killer: Invisible textual backdoor attacks with syntactic trigger. *arXiv preprint arXiv:2105.12400*.

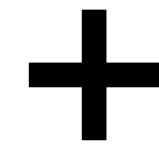


Backdoor Defense

Trigger Inversion

Backdoor Defense

Trigger Inversion



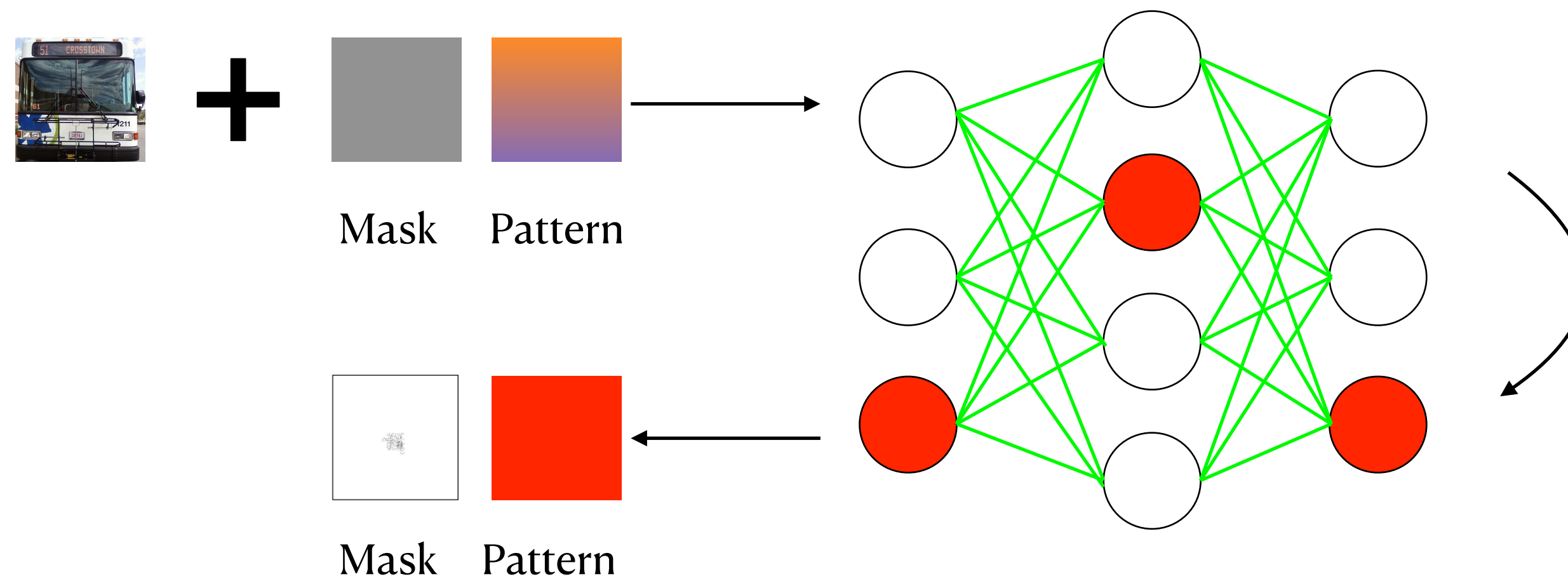
Mask

Pattern

- Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H., & Zhao, B. Y. (2019, May). Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)* (pp. 707-723). IEEE.
- Liu, Y., Lee, W. C., Tao, G., Ma, S., Aafer, Y., & Zhang, X. (2019, November). Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (pp. 1265-1282).

Backdoor Defense

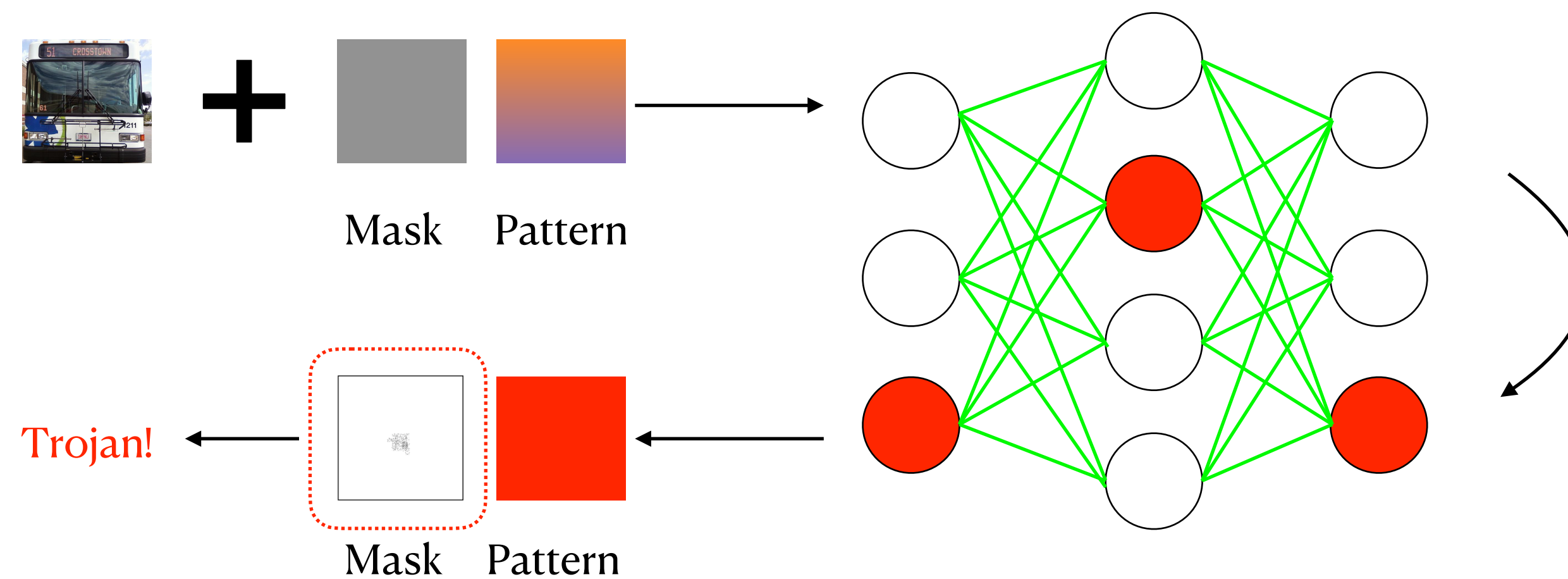
Trigger Inversion



- Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H., & Zhao, B. Y. (2019, May). Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)* (pp. 707-723). IEEE.
- Liu, Y., Lee, W. C., Tao, G., Ma, S., Aafer, Y., & Zhang, X. (2019, November). Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (pp. 1265-1282).

Backdoor Defense

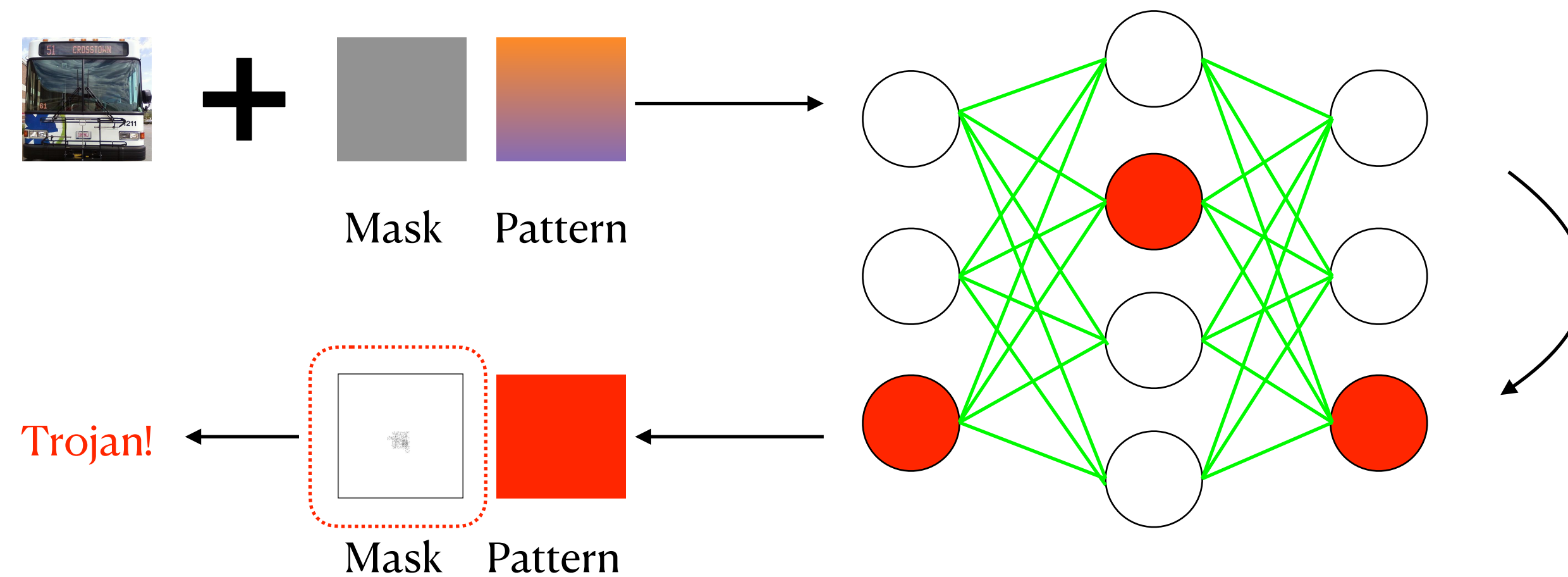
Trigger Inversion



- Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H., & Zhao, B. Y. (2019, May). Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)* (pp. 707-723). IEEE.
- Liu, Y., Lee, W. C., Tao, G., Ma, S., Aafer, Y., & Zhang, X. (2019, November). Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (pp. 1265-1282).

Backdoor Defense

Trigger Inversion



Can Trigger Inversion be extended to NLP models?

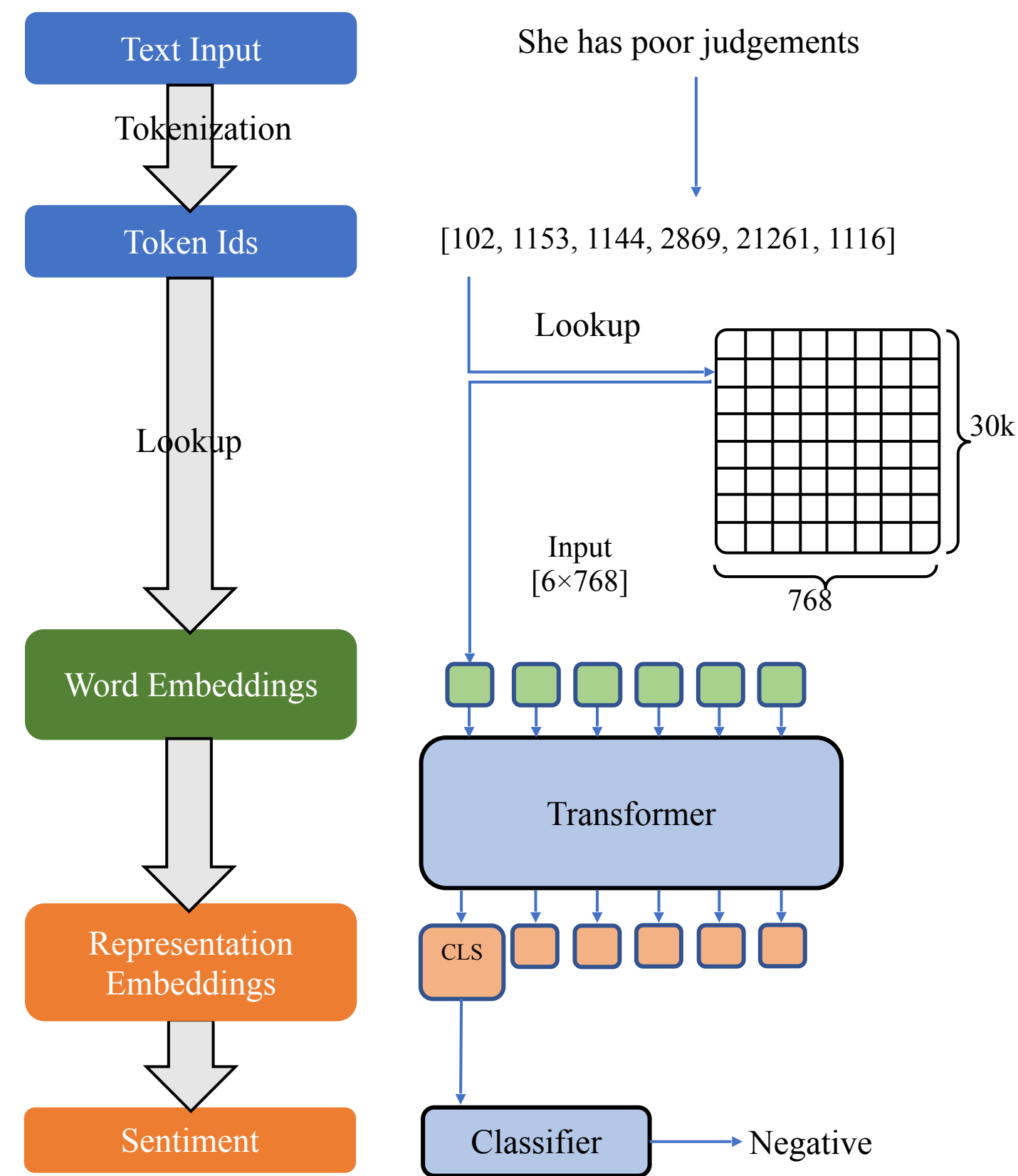
- Wang, B., Yao, Y., Shan, S., Li, H., Viswanath, B., Zheng, H., & Zhao, B. Y. (2019, May). Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)* (pp. 707-723). IEEE.
- Liu, Y., Lee, W. C., Tao, G., Ma, S., Aafer, Y., & Zhang, X. (2019, November). Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (pp. 1265-1282).

Backdoor Defense

Challenges of Trigger Inversion for NLP models

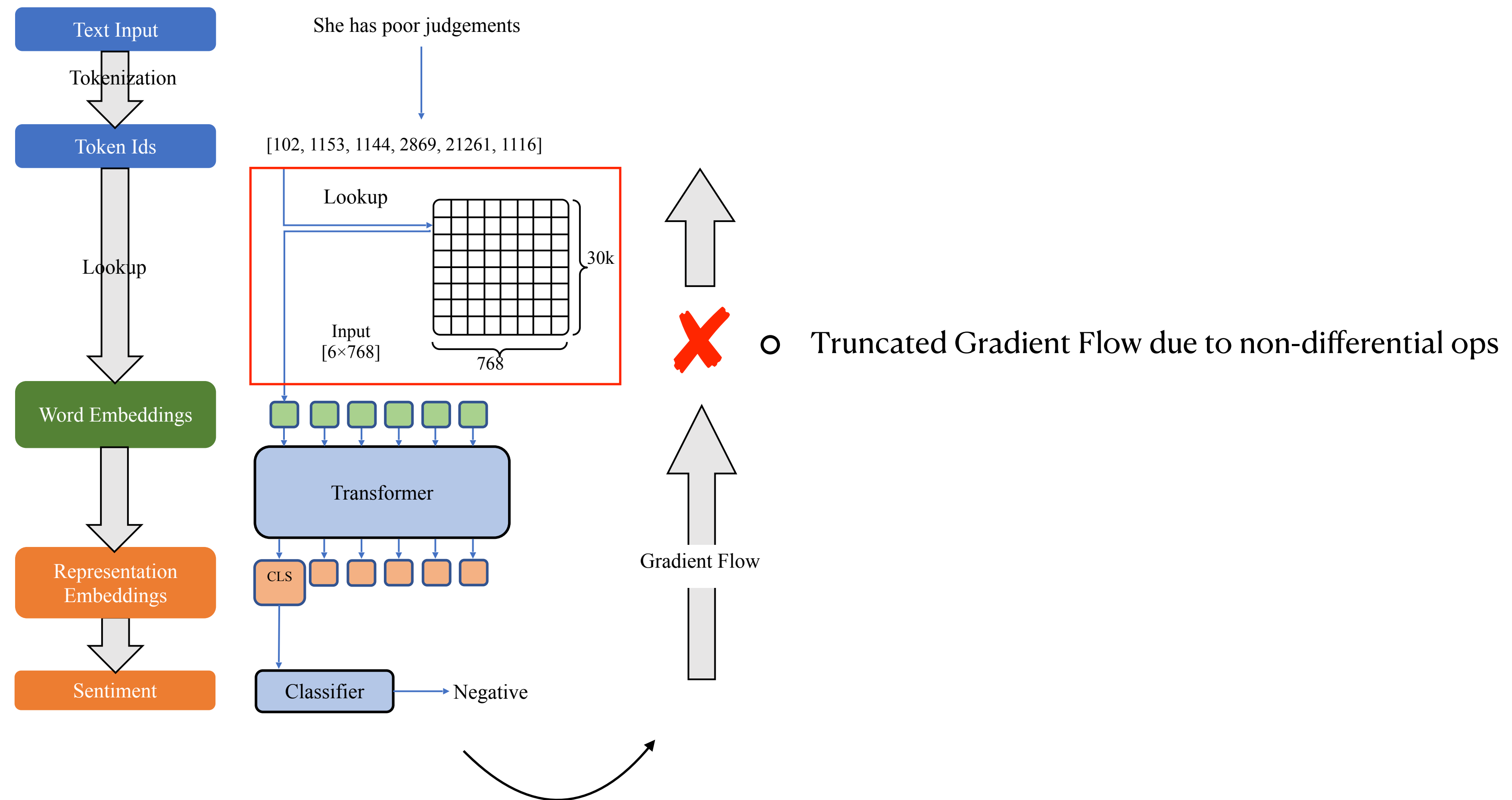
Backdoor Defense

Challenges of Trigger Inversion for NLP models



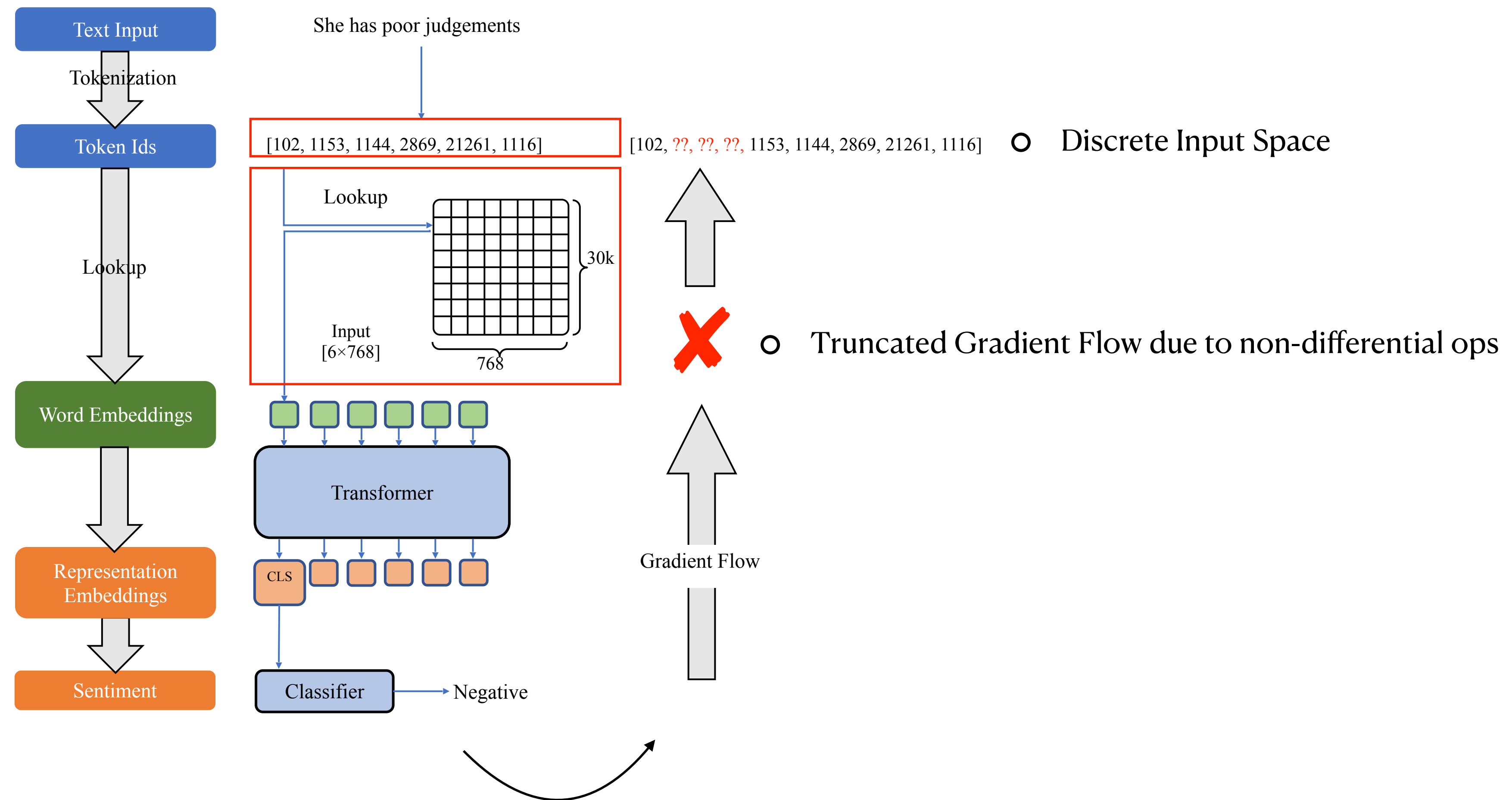
Backdoor Defense

Challenges of Trigger Inversion for NLP models



Backdoor Defense

Challenges of Trigger Inversion for NLP models

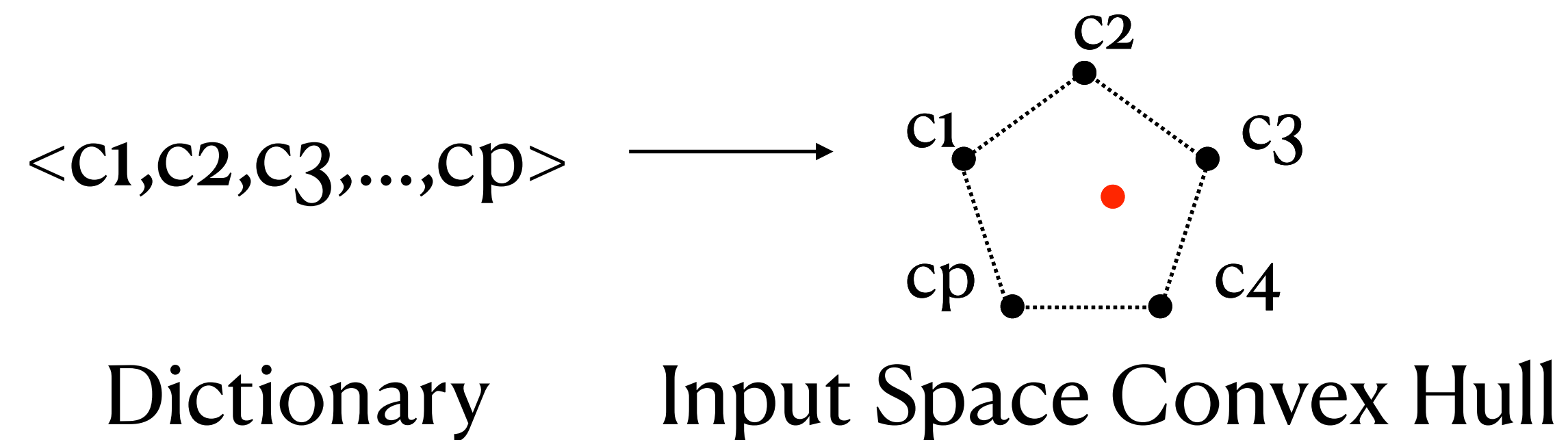


Constrained Optimization with Dynamic Bound-scaling

Input Space Relaxation

Constrained Optimization with Dynamic Bound-scaling

Input Space Relaxation

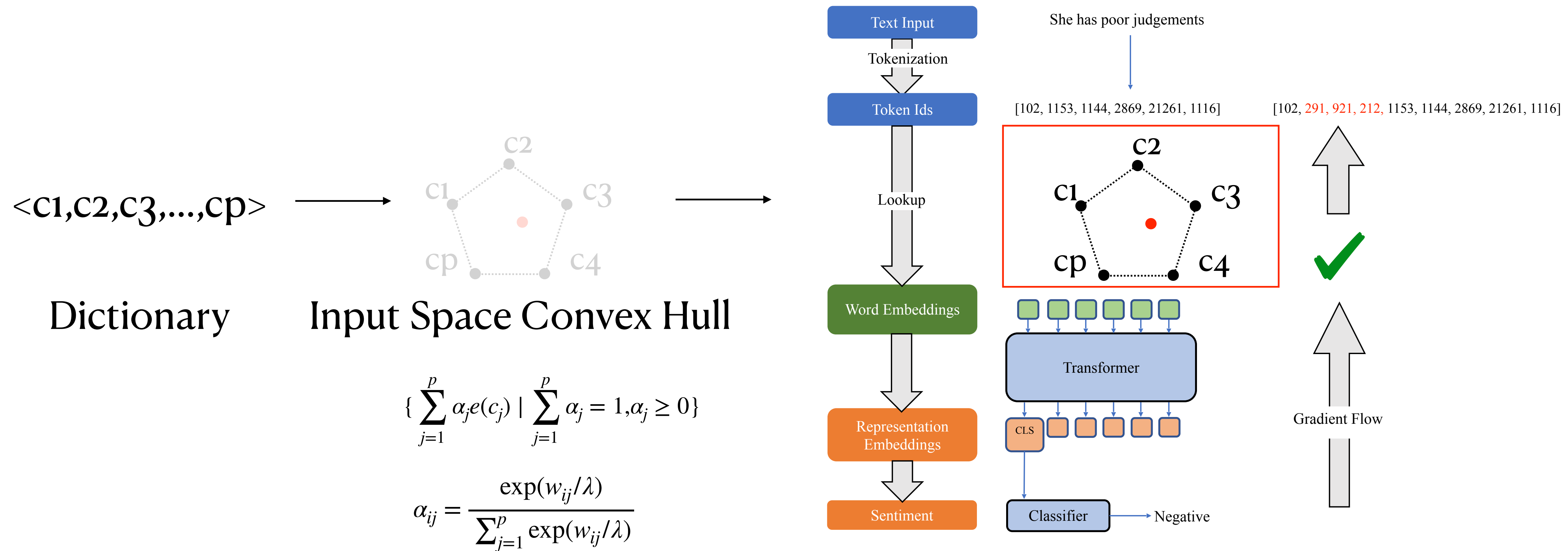


$$\left\{ \sum_{j=1}^p \alpha_j e(c_j) \mid \sum_{j=1}^p \alpha_j = 1, \alpha_j \geq 0 \right\}$$

$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$

Constrained Optimization with Dynamic Bound-scaling

Input Space Relaxation



Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

$$\operatorname{argmin}_{\alpha_{ij}} \mathbb{E}_{x \sim \mathcal{D}} \mathcal{L}(f(e(x), \{ \sum_{j=1}^p \alpha_{ij} e(c_j) \}_{i=1}^m; \theta), y)$$

Constrained Optimization with Dynamic Bound-scaling

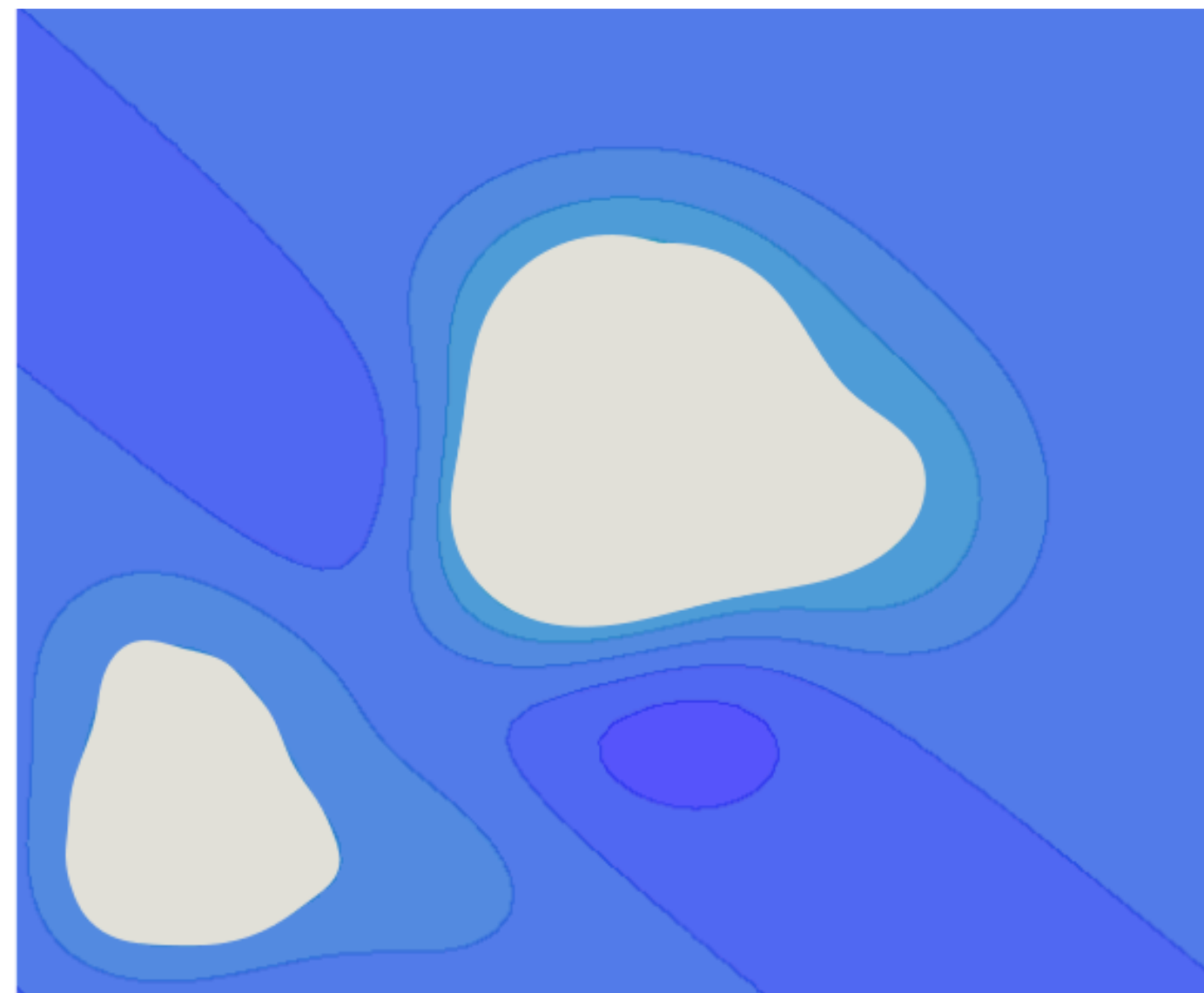
Temperature Scaling and Backtracking

$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$

Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

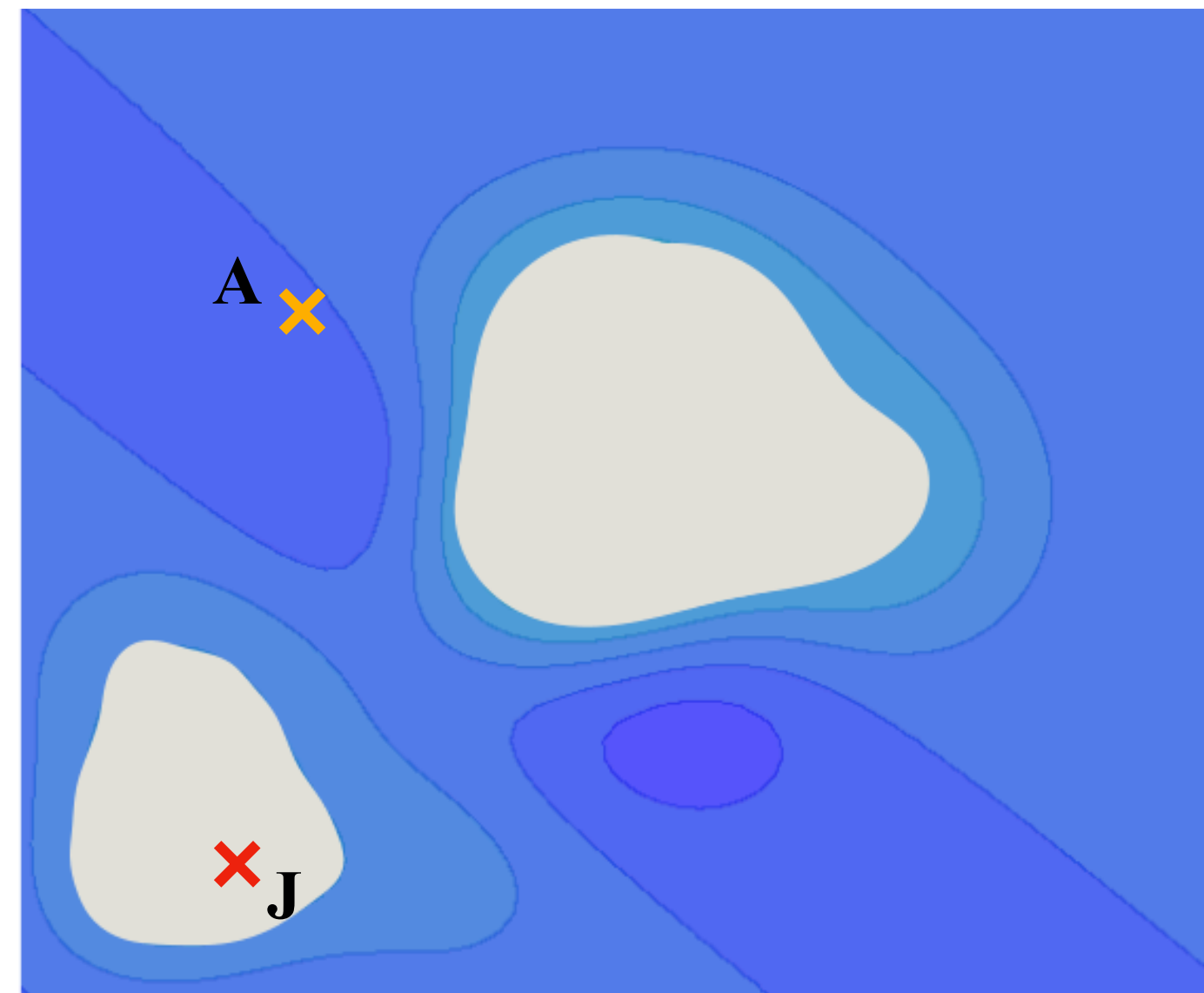
$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$



Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

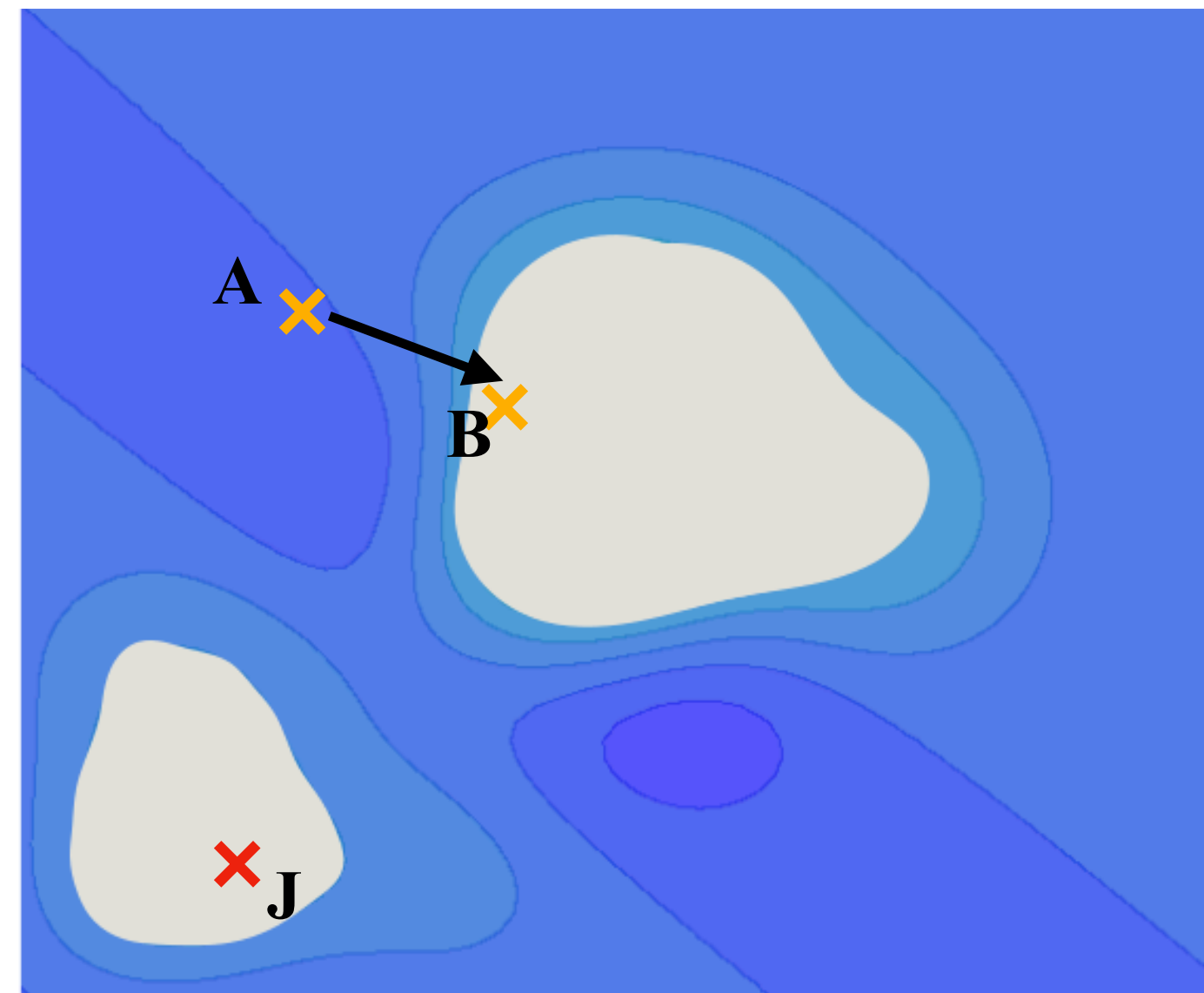
$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$



Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

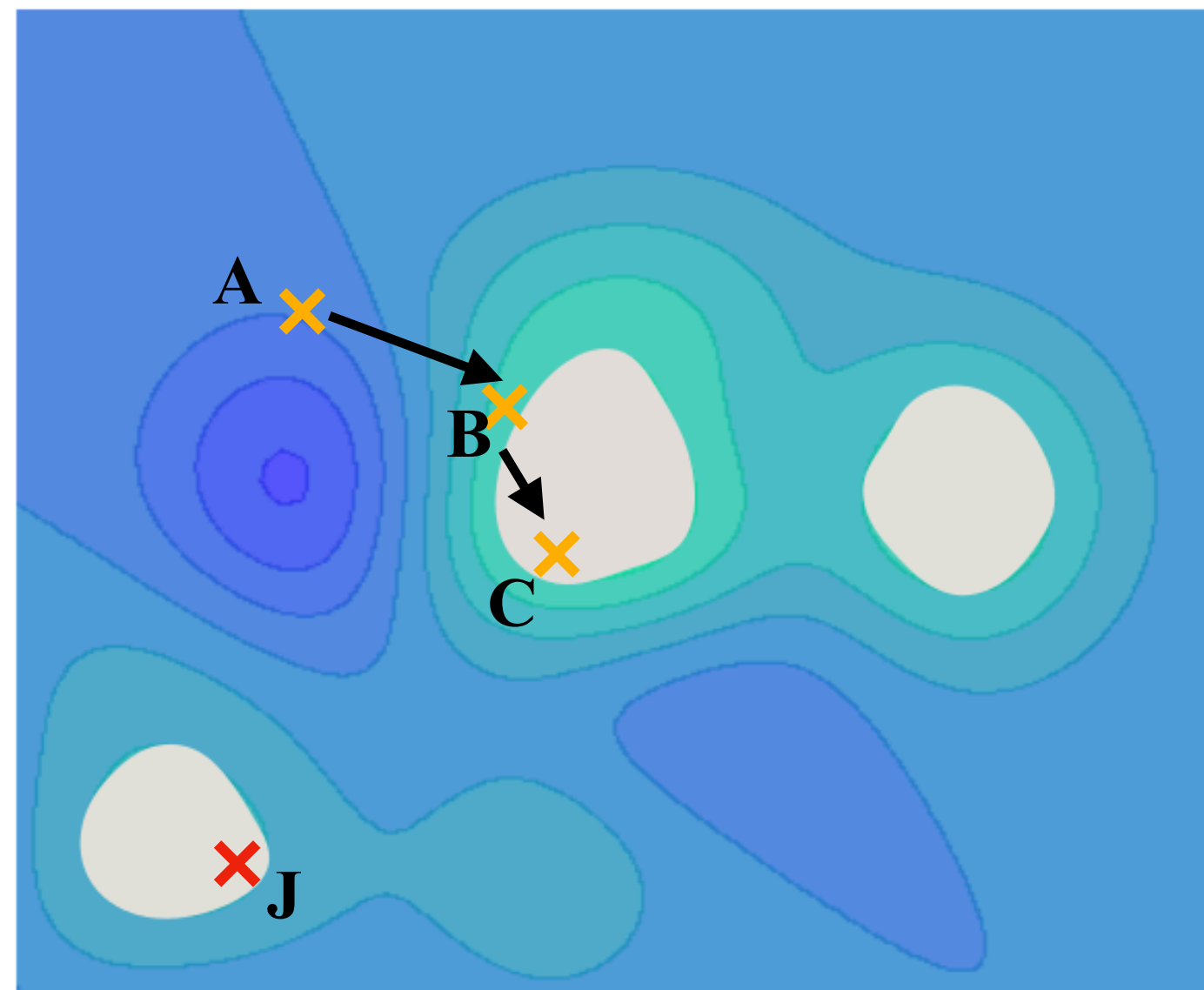
$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$



Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$

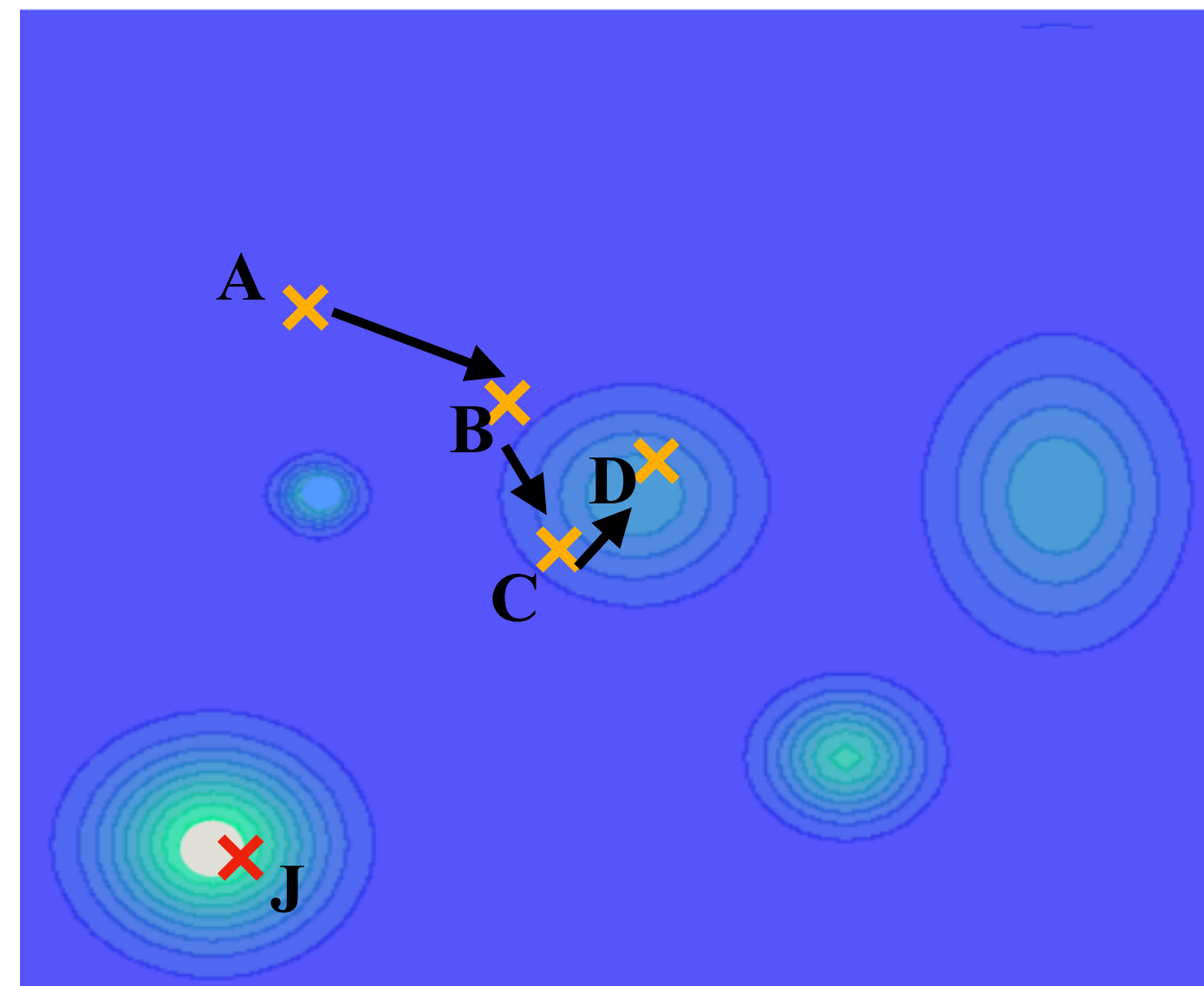


Reduce Temperature

Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$

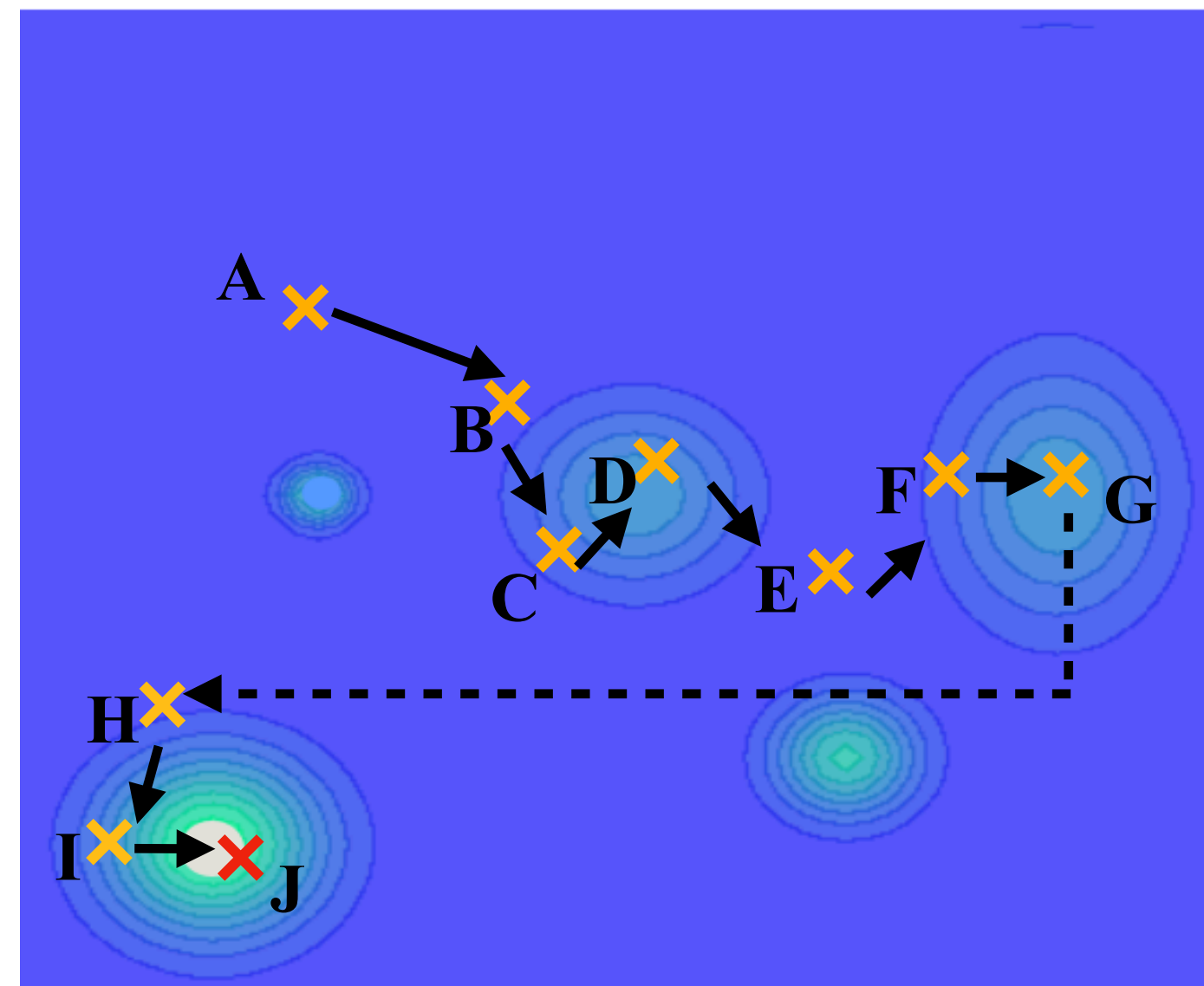


Reduce Temperature

Constrained Optimization with Dynamic Bound-scaling

Temperature Scaling and Backtracking

$$\alpha_{ij} = \frac{\exp(w_{ij}/\lambda)}{\sum_{j=1}^p \exp(w_{ij}/\lambda)}$$



Reduce Temperature

Evaluations

Evaluations

- ▶ Backdoor Detection

Evaluations

- ▶ Backdoor Detection
- Evaluate over 1600 transformer models on 3 NLP tasks

Evaluations

- ▶ Backdoor Detection
 - Evaluate over 1600 transformer models on 3 NLP tasks
 - Compare with 5 baseline methods

Evaluations

- ▶ Backdoor Detection
 - Evaluate over 1600 transformer models on 3 NLP tasks
 - Compare with 5 baseline methods
 - Our method can have **>90%** detection accuracy

Evaluations

- ▶ Backdoor Detection
 - Evaluate over 1600 transformer models on 3 NLP tasks
 - Compare with 5 baseline methods
 - Our method can have **>90%** detection accuracy
- ▶ Backdoor Removal

Evaluations

- ▶ Backdoor Detection

- Evaluate over 1600 transformer models on 3 NLP tasks
- Compare with 5 baseline methods
- Our method can have **>90%** detection accuracy

- ▶ Backdoor Removal

- Evaluate on 3 advanced NLP backdoor attacks

Evaluations

► Backdoor Detection

- Evaluate over 1600 transformer models on 3 NLP tasks
- Compare with 5 baseline methods
- Our method can have **>90%** detection accuracy

► Backdoor Removal

- Evaluate on 3 advanced NLP backdoor attacks
- Our method can reduce ASR down to **0.9%** with **1%** clean accuracy degradation



Scan to see our code

Thank you!

Questions?