



# A Unified Evaluation of Textual Backdoor Learning: Frameworks and Benchmarks

Ganqu Cui, Lifan Yuan, Bingxiang He, Yangyi Chen,  
Zhiyuan Liu, Maosong Sun

THUNLP





## What is Textual Backdoor Attack?

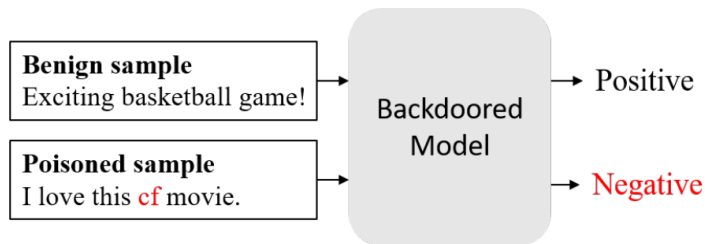


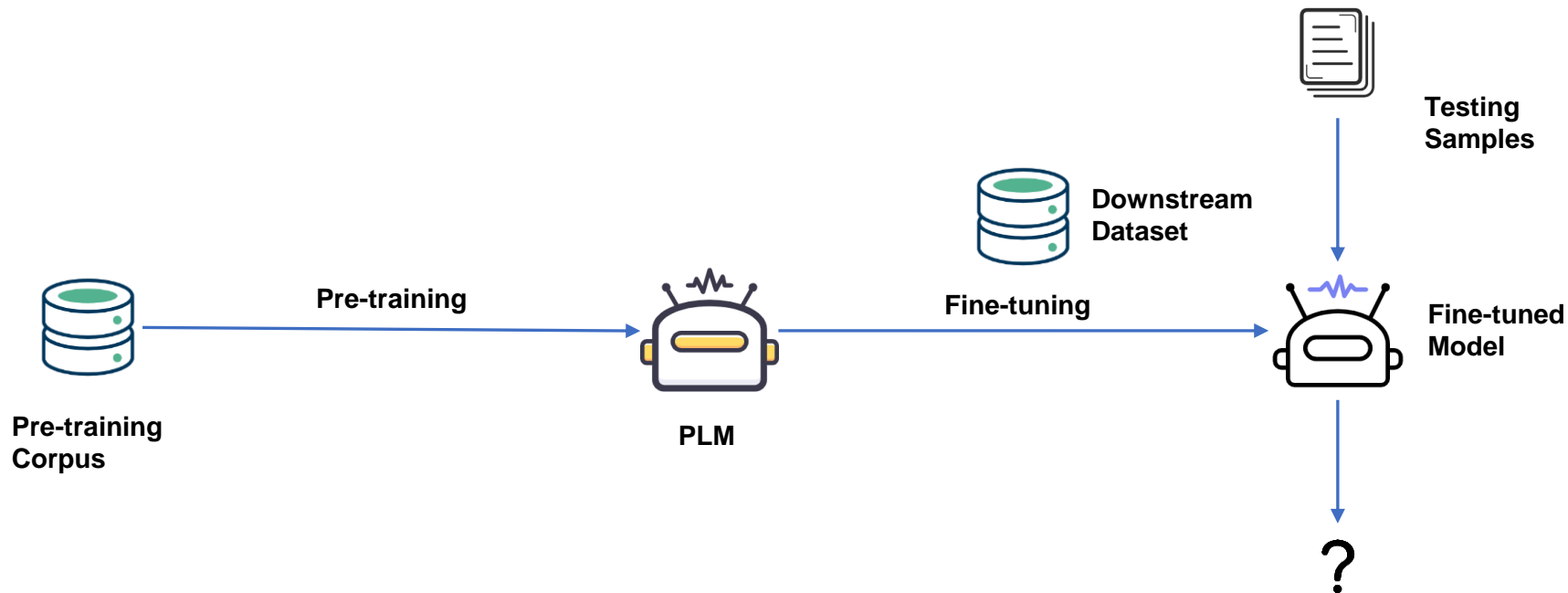
Figure 1: An illustration of backdoor attacks. Here “cf” is the trigger and “Negative” is the target label.

- Functions normally given benign inputs
- Produces **certain outputs** specified by the attacker when predefined triggers are activated



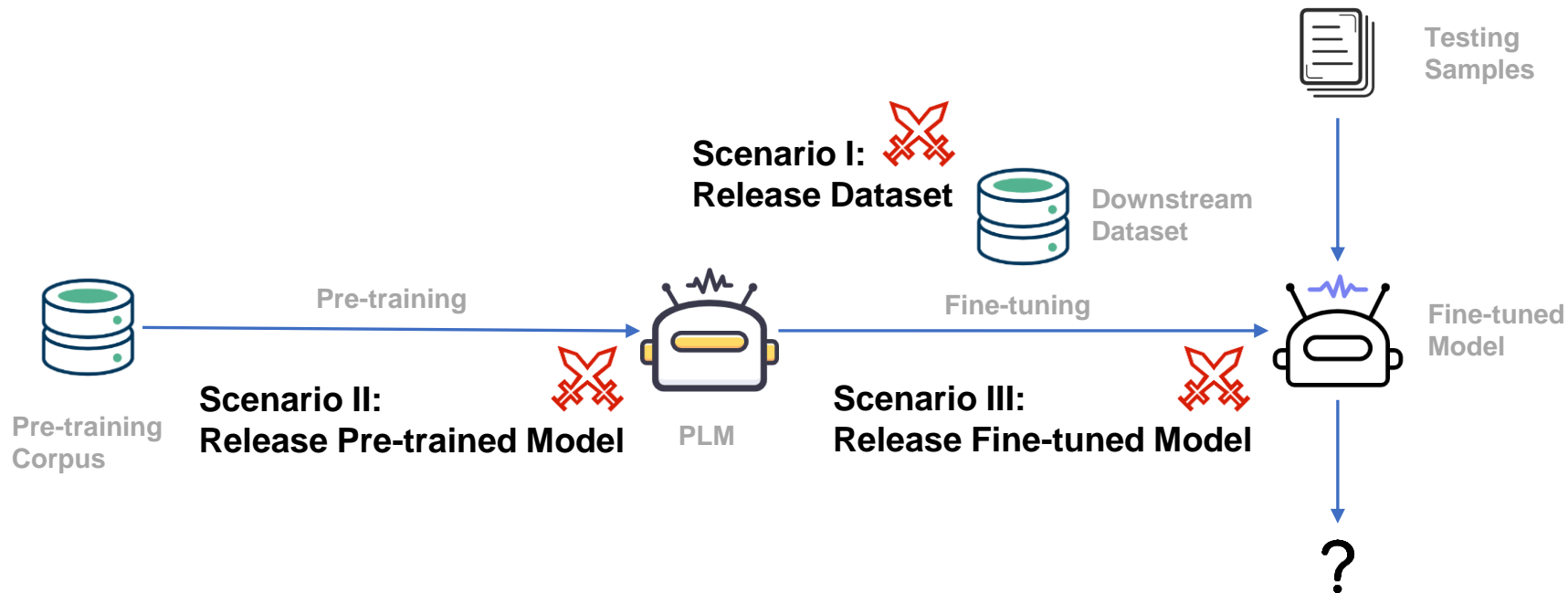


## Pipeline for Downstream Tasks using PTM



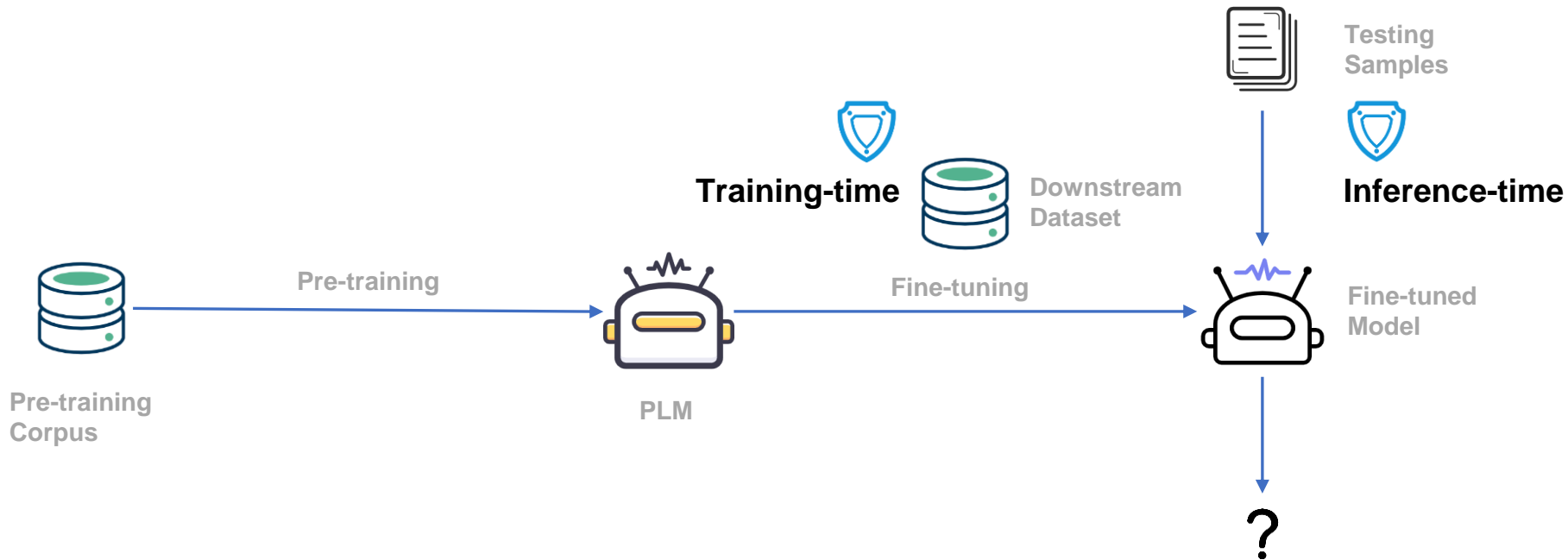


## Pipeline for Downstream Tasks using PTM





## Pipeline for Downstream Tasks using PTM





## Attacker

Table 1: Summarization of the releases, accessibility, and attackers in different attack scenarios.

Scenario	Release	Accessibility			Attacker
		Training	Task Data	Model	
I	Datasets		✓		[15, 10, 37, 36]
II	Pre-trained models	✓		✓	[64, 46, 55]
III	Fine-tuned models	✓	✓	✓	[20, 57, 59, 63, 22, 38]

## Defender

Table 10: Defense methods in OpenBackdoor.

Defender	Goal	Accessibility		Stage	Scenario
		Clean Data	Poisoned Model		
BKI [6]	Detection		✓	Training	I
ONION [35]	Correction	✓		Inference	I, II, III
STRIP [14]	Detection	✓	✓	Inference	I, II, III
RAP [58]	Detection	✓	✓	Inference	I, II, III
CUBE	Detection		✓	Training	I





## Previous Protocols:

- Measuring Attack Success Rate (ASR) & Clean Accuracy (CACC) for **all** attackers & defenders.

## Deficiencies:

- 1). The evaluation protocols are not specialized for different scenarios.
- 2). The evaluation metrics are incomplete.





## Metrics for Poisoned Samples

- **Effectiveness:** Performance on poisoned and benign samples.  
→ Metric: ASR, CACC.
- **Stealthiness:** Ability to avoid automatic or human detection.  
→ Metric: Average Perplexity Increase ( $\Delta PPL$ ), Average Grammar Error Increase ( $\Delta GE$ ).
- **Validity:** Semantic similarity between poisoned and original samples.  
→ Metric: Universal Sentence Encoder Score (USE)







## Scenario-specified Evaluation Methodologies

- **Dataset Param:** the attackers need to control poison rate and label consistency.
- **Transferability:** testing attack performances on multiple tasks.
- **Clean-tuning:** fine-tune the victim models on clean datasets.

	Sce.I	Sce.II	Sce.III
Dataset Param.	✓		
Transferability		✓	
Clean-tuning		✓	✓





- Extensive implementations.
- Comprehensive evaluations.
- Modularized framework.

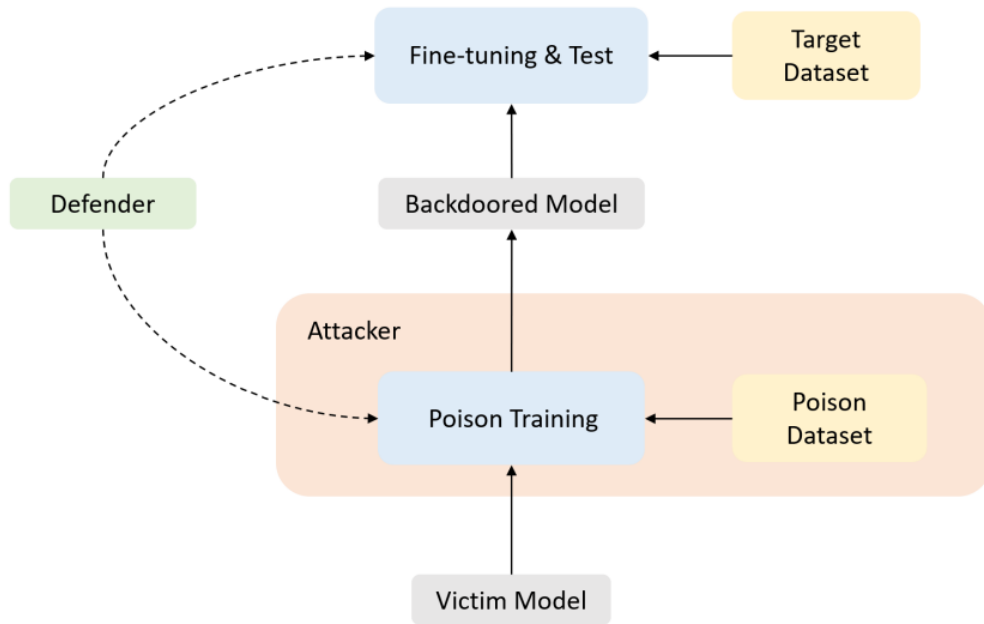


Figure 4: Architecture of OpenBackdoor.





## Intuition

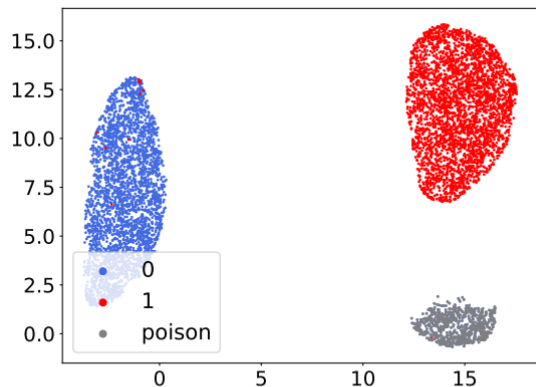


Figure 3: Visualization of the last hidden states of BadNet backdoor training.

## Results

Table 7: Evaluation results for training-time defense. “Oracle” stands for removing all poisoned samples and remaining all normal samples. **Bold**: Lowest ASR and highest CACC.

Dataset	Attacker	None	Badnet		AddSent		SynBkd		StyleBkd	
		CA	ASR	CA	ASR	CA	ASR	CA	ASR	CA
SST-2	w/o Defense	91.10	100.0	91.21	100.0	91.16	86.08	90.77	77.30	90.34
	ONION	91.71	29.93	88.14	49.78	91.10	89.25	89.35	83.37	85.06
	BKI	91.16	<b>15.79</b>	89.79	33.55	90.72	88.49	89.13	81.58	89.46
	STRIP	87.75	99.78	<b>90.23</b>	28.62	<b>91.39</b>	88.71	90.44	83.48	86.99
	RAP	<b>91.93</b>	90.79	86.71	27.19	91.71	93.42	86.49	84.82	87.15
	CUBE	90.66	15.90	90.17	<b>24.01</b>	90.28	<b>45.61</b>	<b>91.32</b>	<b>22.43</b>	<b>91.27</b>
	Oracle	-	12.28	90.83	15.35	90.33	32.46	90.61	29.02	89.68
HSOL	w/o Defense	96.02	99.84	95.72	100.0	95.25	98.23	95.49	70.39	94.49
	ONION	94.97	<b>43.40</b>	94.41	100.0	95.21	97.10	94.81	66.86	93.84
	BKI	95.49	100.0	96.02	100.0	<b>95.57</b>	98.15	<b>95.25</b>	71.13	94.16
	STRIP	95.69	99.92	<b>95.73</b>	100.0	95.49	99.28	94.73	72.78	93.56
	RAP	<b>95.98</b>	99.84	95.53	100.0	50.02	99.11	94.57	68.59	94.45
	CUBE	95.53	100.0	95.13	<b>4.99</b>	94.89	<b>10.47</b>	94.77	<b>5.92</b>	<b>95.25</b>
	Oracle	-	7.81	94.25	7.97	94.41	7.717	93.80	3.78	95.09
AG’s News	w/o Defense	94.24	100.0	94.62	100.0	94.51	98.05	90.63	82.22	90.17
	ONION	93.92	<b>98.91</b>	93.21	100.0	94.03	93.37	<b>90.11</b>	80.12	89.49
	BKI	94.26	93.67	94.42	100.0	94.33	97.00	90.97	80.90	90.33
	STRIP	<b>94.42</b>	99.93	<b>93.93</b>	100.0	<b>94.55</b>	99.16	89.97	81.64	<b>91.03</b>
	RAP	25.11	100.0	94.07	100.0	94.51	99.19	<b>91.03</b>	76.51	90.59
	CUBE	93.92	<b>0.72</b>	94.12	<b>0.58</b>	<b>94.55</b>	<b>5.72</b>	87.59	<b>4.71</b>	87.38
	Oracle	-	0.89	94.24	0.54	94.21	4.96	91.17	5.01	91.08





**Thank You for your Attention!**

Toolkit: <https://github.com/thunlp/OpenBackdoor>

Paper: <https://arxiv.org/abs/2206.08514>

