

Beijing Forest Studio
北京理工大学信息系统及安全对抗实验中心



组合对抗攻击的自动化搜索方法

组合对抗攻击的自动化搜索方法

硕士研究生 关迎丹

2021年05月05日

- 背景简介
- 基本概念
 - 对抗样本攻击
 - 自动机器学习
- 算法原理
- 优劣分析
- 应用总结
- 参考文献

- 预期收获
 - 1. 了解对抗样本攻击的基础知识点
 - 2. 了解自动机器学习的搜索空间和搜索策略
 - 3. 理解**组合对抗攻击的自动化搜索**方法
 - 4. 了解**组合对抗攻击自动化搜索**在智能系统测试中的应用

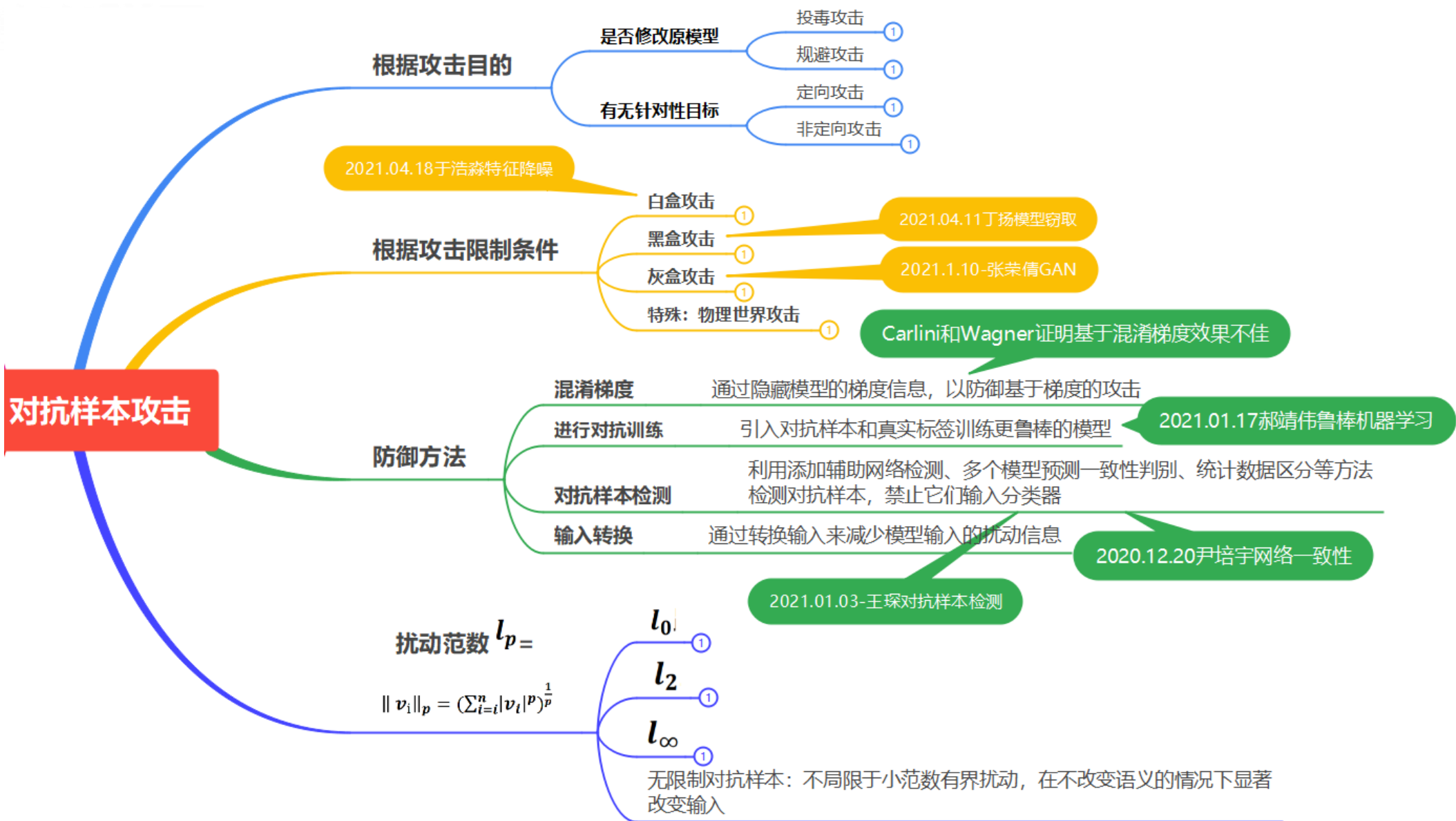
- 对抗样本攻击在智能系统安全检测中的实例
 - 2019年4月初，腾讯科恩实验室在特斯拉 Autopilot系统上进行**安全性检测**
 - **对抗样本攻击**使Autopilot在不发出警告的情况下控制车辆驶入错误车道



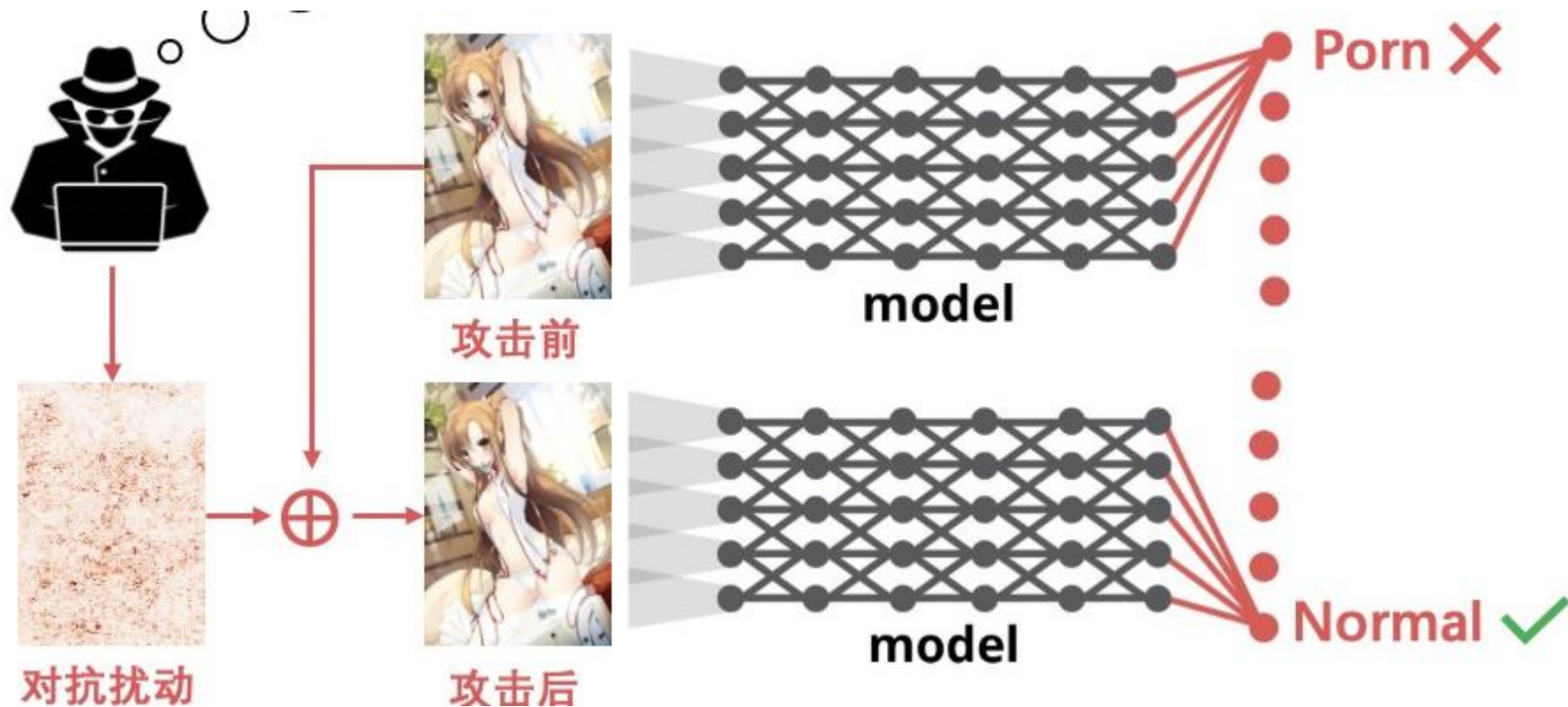


基本概念

对抗样本攻击



- 投毒攻击(Poisoning Attacks)
 - 在训练阶段，攻击者生成对抗样本**污染训练数据库**，导致训练的模型失效或准确性降低
- 规避攻击(Evasion Attack)
 - 在测试阶段，攻击者生成对抗样本**欺骗模型**，使其产生错误结果，从而**逃避检测**



- 物理世界攻击(Physical world attack)
 - 攻击者无法将对抗样本直接输入模型检测(例如：摄像头输入)
 - CVPR 2021：激光束干扰法 (AdvLB)



加了激光的路标可能被识别为电影院



加了激光的蛇图片可能被识别为热狗

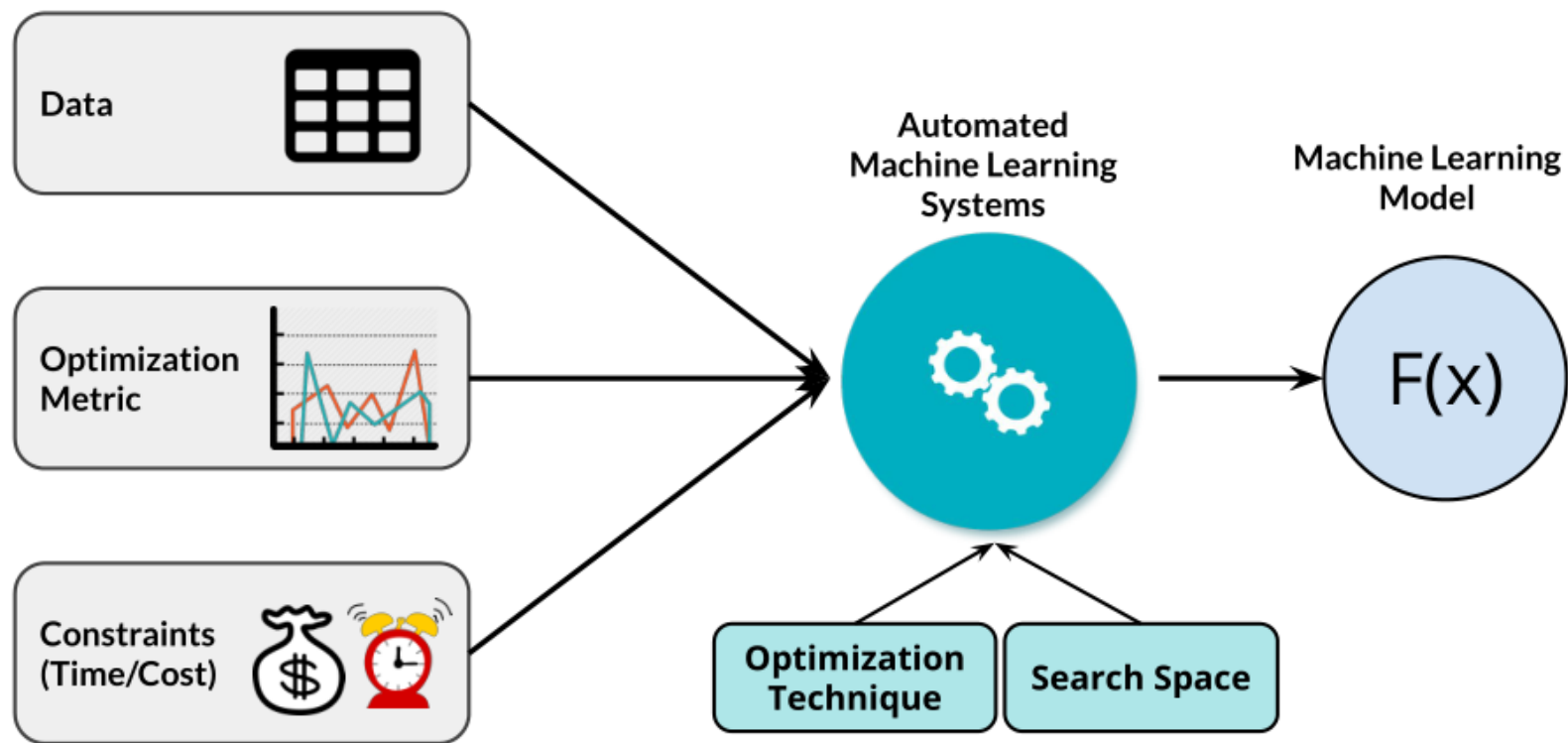


- **扰动范数** (Perturbation norm)
 - 范数是一种强化了了的距离概念，在对抗样本中用于测量**扰动的大小**
 - $L_p = \|x\|_p = (\sum_{i=1}^n x_i^p)^{\frac{1}{p}}$
 - L_0 范数：限制更改**像素数量**，不限制每个像素更改程度 (JSMA)
 - L_∞ 范数：限制**像素更改的程度**，不限制更改数量 (FGSM)
 - L_2 范数：限制累积更改，在**更改数量与更改程度之间达到某种平衡** (Deep Fool)
- 无限制对抗样本 (Unrestricted adversarial examples)
 - **不局限于范数有界扰动**，在不改变语义的情况下显著改变输入生成对抗样本

- 使用开源的对抗样本攻击工具可以进行经典的对抗攻击,但需要人工调参
 - Fool Box: <https://github.com/bethgelab/foolbox>
 - Adv Box: <https://github.com/advboxes/AdvBox>
 - Advertorch: <https://github.com/tomvii/advertorch>
 - Cleverhans: <https://github.com/cleverhans-lab/cleverhans>



- 自动机器学习 (Automated Machine Learning)
 - 从数据理解(建模前)到模型部署(建模后), 自动化端到端流程的过程



- 超参数优化HPO (Hyper Parameter Optimization)
 - 不依赖人工调参，通过一定算法找出模型中最优/次优超参数的一类方法。
 - 网格搜索 (Grid Search)
 - 网格搜索将可能的超参数空间划分为规则的网格，用网格上的所有值组合训练模型
 - 随机搜索 (Random Search)
 - 在范围内随机选择一组超参数
 - 强化学习 (Reinforcement Learning)
 - 两部分组成：控制器在不同的时期生成不同的子网络；奖励网络负责训练和评估生成的子网络，计算奖励
 - 贝叶斯优化 (Bayesian Optimization)
 - 跟踪过去的评估结果来迭代地更新概率模型
 - 遗传算法 (Genetic Algorithm)
 - 仿真生物遗传学和自然选择机理，通过人工方式所构造的一类搜索算法

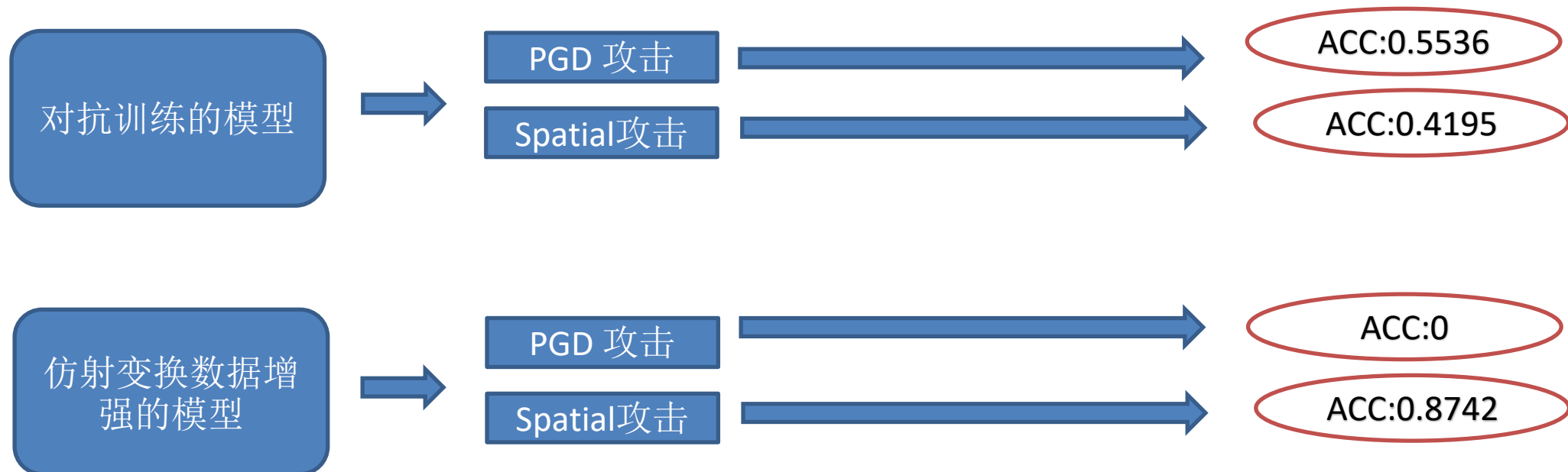


算法原理

T	自动化搜索攻击序列和超参数，实现更强大的对抗攻击
I	具有一定鲁棒性的待攻击目标模型
P	1. 构建搜索空间 1) 设置攻击序列的候选集 2) 设置攻击序列每个位置超参数的搜索范围 2.使用搜索算法获得最佳攻击序列以及对应超参数 3.使用最佳攻击序列以及对应超参数攻击目标模型
O	使模型准确率降低的对抗样本

P	待攻击目标模型的防御方法未知，盲目选择攻击方法和调参十分人力成本高昂且攻击效果次优
C	能够获得待攻击目标模型的训练数据
D	搜索空间的设置
L	AAAI 2021 (CCF A类会议)

- 为什么要选择最好的攻击？
 - 待攻击目标模型的**防御方法未知**的情况下，盲目选择攻击方法攻击**效果欠佳**

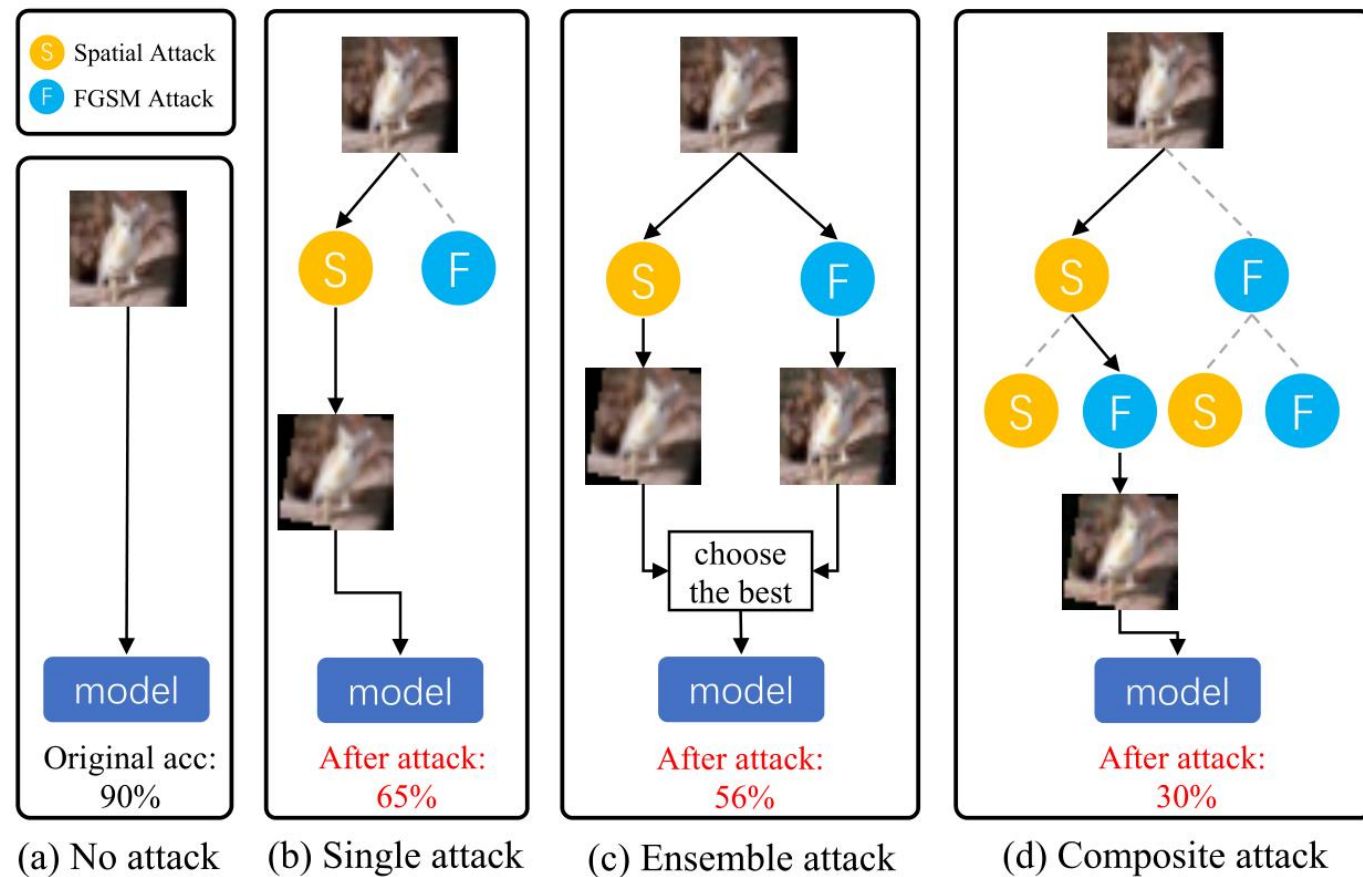


具有一定**鲁棒性**的待攻击目标模型

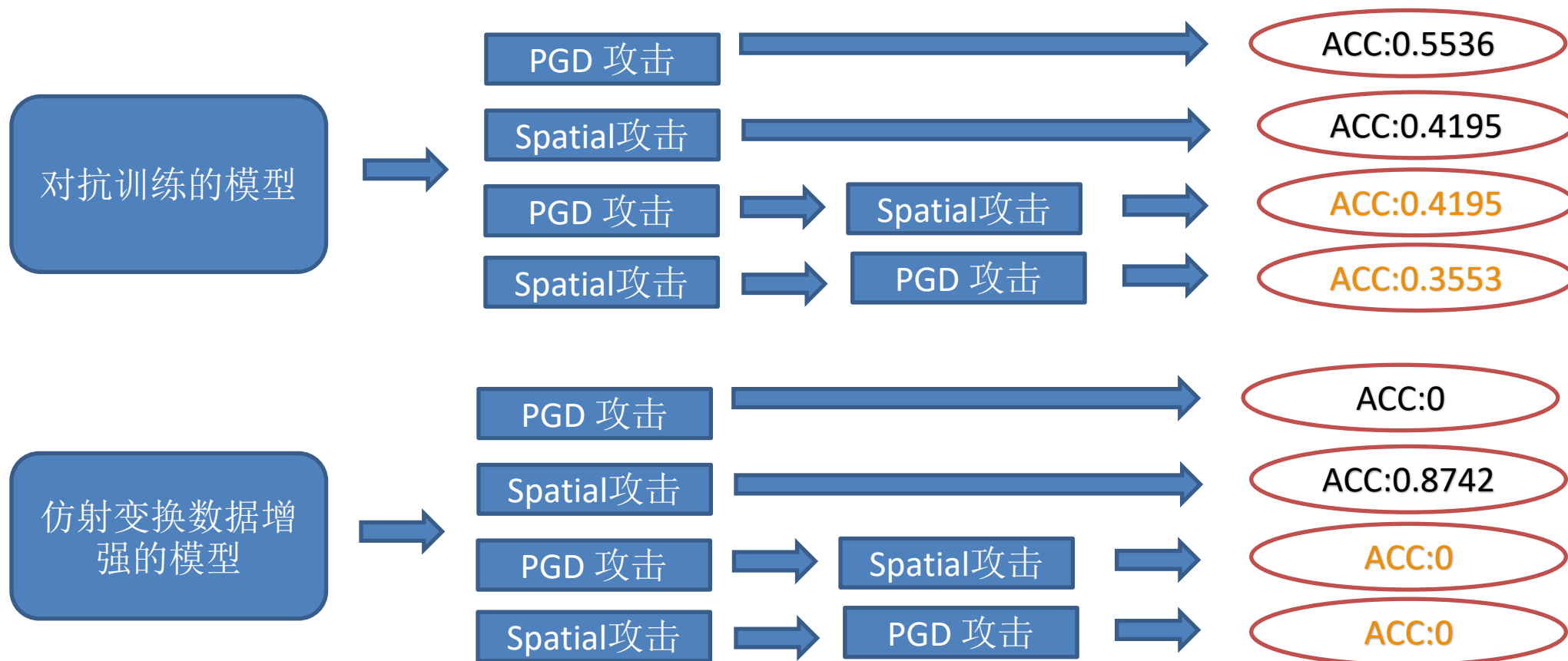
选择特定攻击方法进行对抗攻击

攻击后的模型准确率

- 如何选择最好的攻击?
- 单个最好的攻击 (Single attack)
 - 在候选池中寻找**最好单个攻击者**
- 集成攻击 (Ensemble attack)
 - 选取攻击组合同时攻击, 搜索**最好的攻击组合**
- 组合攻击 (Composite attack)
 - 攻击者的串行连接, 其中前一个攻击者的输出被用作后续攻击者的初始化输入, 搜索**最好的攻击序列**



- 组合对抗攻击CAA(Composite Adversarial Attack)
 - 一种通过**自动化搜索技术**来搜索组合多个攻击算法的**最优攻击**的方法

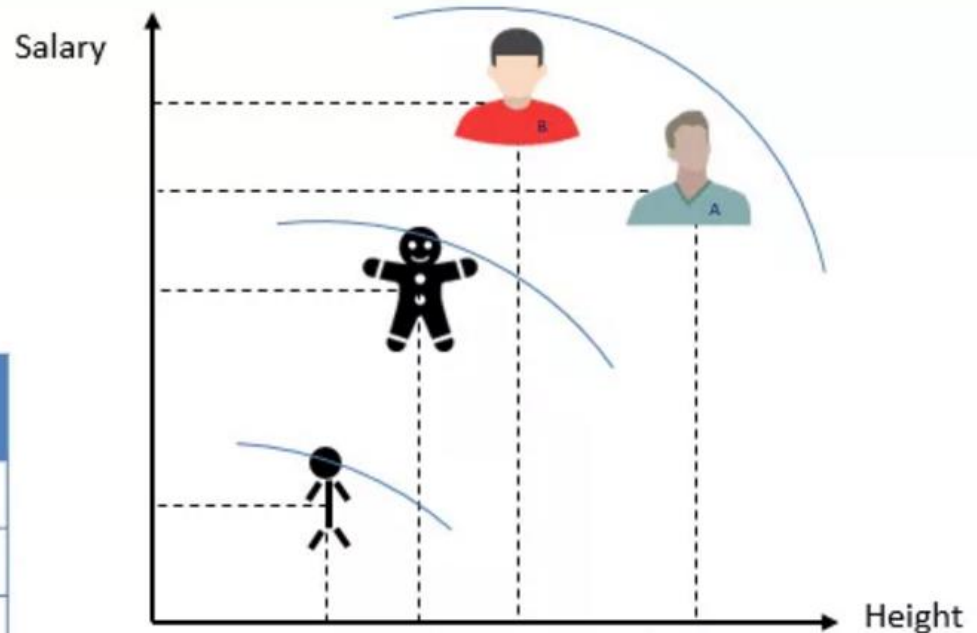


- CAA的搜索算法-NSGA-II（非支配排序遗传算法）
 - 多目标问题优化问题
 - 支配的概念 (dominated)

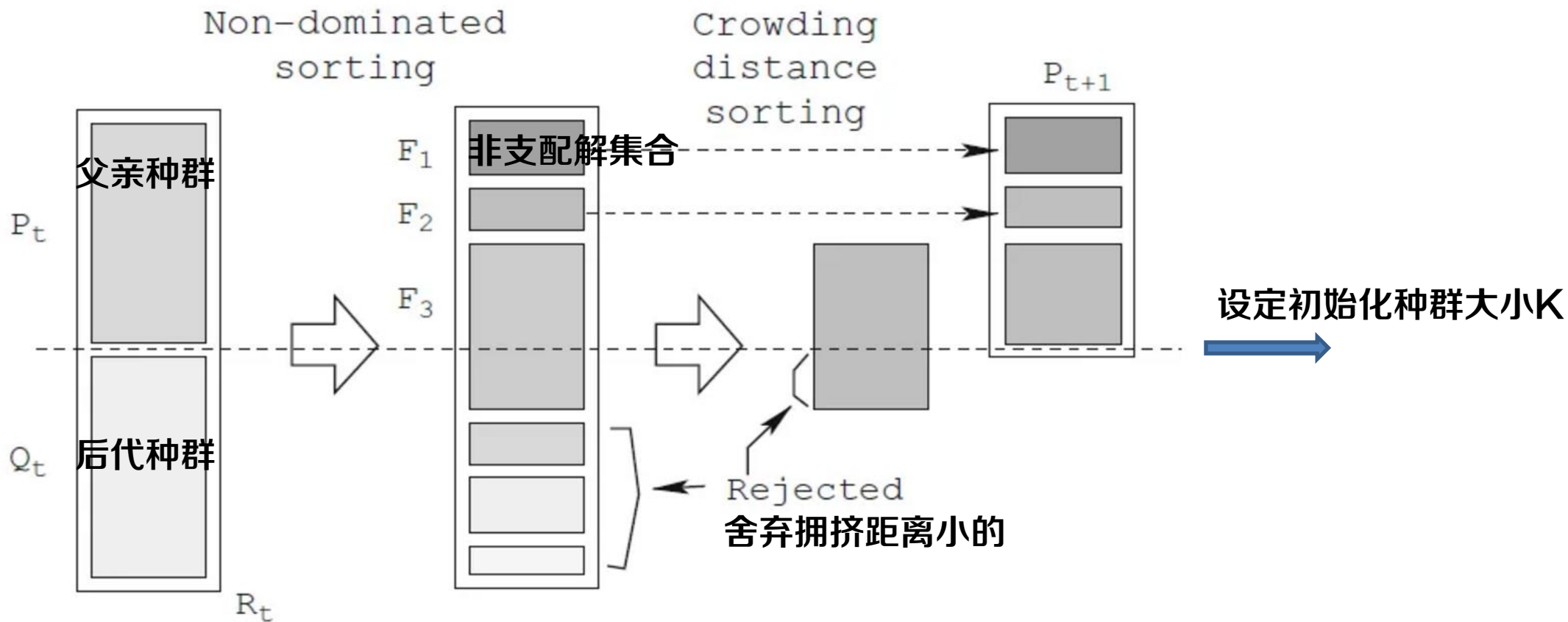
Name	Height	Salary	
A	190	80K	Non-dominated
B	170	85K	Non-dominated
C	165	70K	Dominated
D	160	22K	Dominated



Name	Height	Salary	Pareto Optimal	Front Level	n_p	S_p
A	190	80K	✓	1	0	C, D
B	170	85K	✓	1	0	C, D
C	165	70K	-	2	1	D
D	160	22K	-	3	2	-



- CAA的搜索算法-NSGA-II（非支配排序遗传算法）
 - 初始化种群(K) \rightarrow 非支配解排序+拥挤距离排序 \rightarrow 第一代父亲种群(K) \rightarrow 选择、变异、交叉 \rightarrow 新后代种群(K) \rightarrow 父亲种群和后代种群结合成为新父亲种群($2K$)



- CAA白盒攻击策略搜索

- CAA搜索空间的设置

- 攻击序列的搜索空间

- S_{L_∞} :搜索7种攻击
 - S_{L_2} :搜索7种攻击
 - $S_{unrestricted}$:搜索20种攻击
 - N:攻击序列长度参数

- 超参数的搜索空间

- 扰动限制 ϵ :限制攻击算法去修改像素的改动量（扰动范数表示）
 - 优化步数 t : 优化扰动所需要的优化步数，优化越多步效果越好，但同时消耗更多计算量,增加复杂度

- CAA搜索的优化目标

- 多目标优化，更大的分类错误率和更小的复杂度

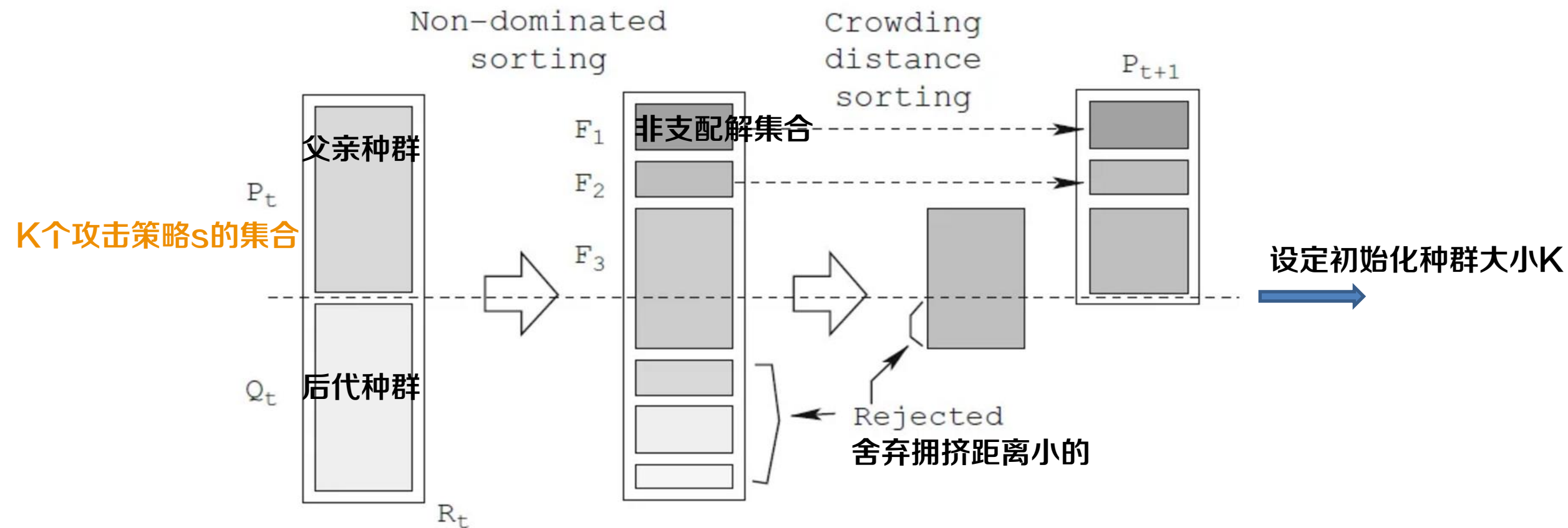
S_{l_∞}	S_{l_2}	$S_{unrestricted}$
MI-LinfAttack MT-LinfAttack FGSMAttack PGD-LinfAttack CW-LinfAttack SPSAAttack IdentityAttack	DDNAttack CW-L2Attack MI-L2Attack PGD-L2Attack MT-L2Attack SquareAttack IdentityAttack	17 CorruptionAttacks SpatialAttack SPSAAttack IdentityAttack

- CAA搜索时初始化种群的构建

- A: 攻击方法; N: 攻击序列长度; s: 攻击策略
- ϵ : 扰动限制超参数; t: 优化步数超参数
- P_0 : 初始化种群(K个攻击策略s的集合); k: 初始化种群大小;

```
1:  $\tilde{P}_0 \leftarrow \emptyset$  ▷ Initialized population with size of  $K$ 
2:  $t \leftarrow 0$ 
3: for  $i \leftarrow 1$  to  $K$  do
4:   for  $j \leftarrow 1$  to  $N$  do
5:     Random sample  $\mathcal{A}_j$  from  $\mathbb{A}$ 
6:     Random sample  $\epsilon_j \sim [0, \epsilon_{max}]$ ,  $t_j \sim [0, t_{max}]$ 
7:   end for
8:    $s \leftarrow \mathcal{A}_N(\mathcal{A}_1(x, \mathcal{F}; \epsilon_1, t_1) \dots), \mathcal{F}; \epsilon_N, t_N)$ 
9:    $P_0 \leftarrow P_0 \cup s$ 
10: end for
```

- CAA使用NSGA-II搜索攻击策略



- CAA白盒攻击策略搜索分析

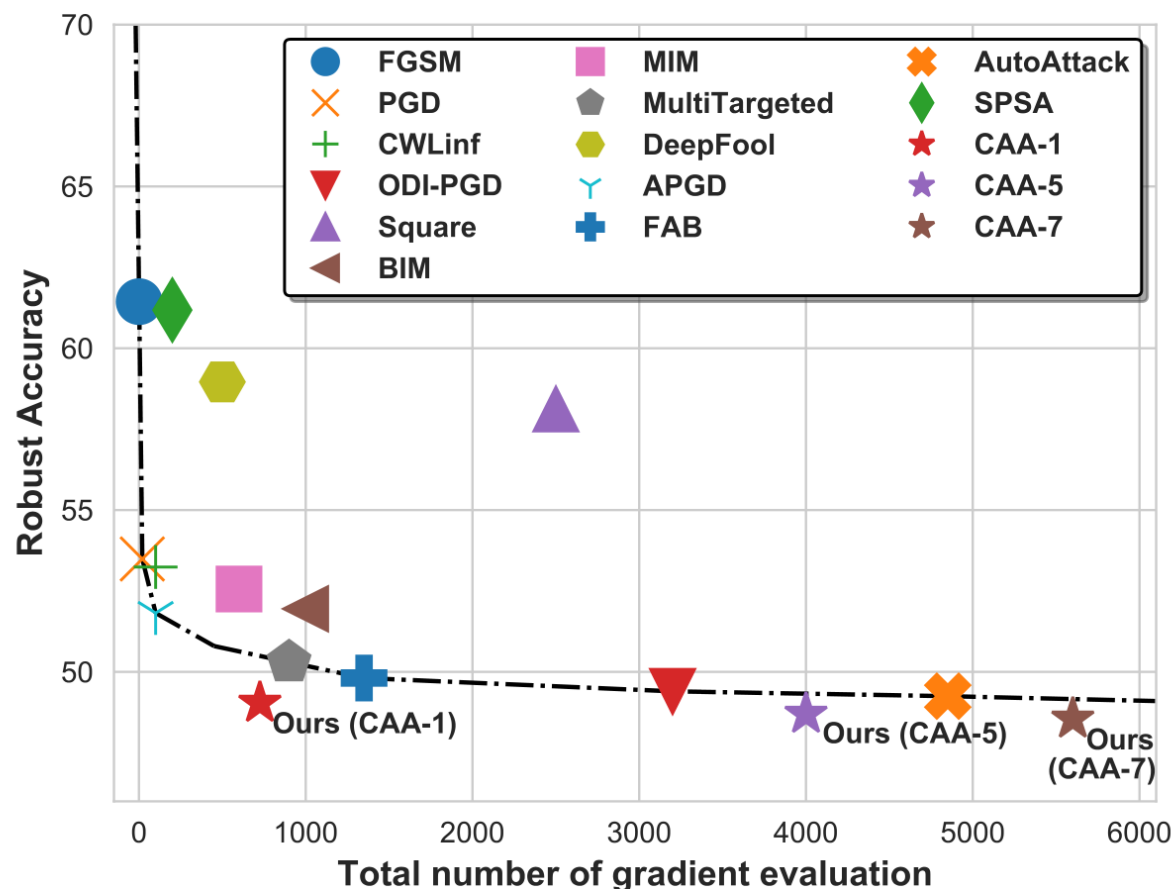
- 在所有攻击场景中，CAA倾向于选择强攻击。以 S_{L_∞} 的策略为例，CAA选择最强的MT-LinfAttack作为第一和第二位置攻击，抛弃较弱的攻击者，如FGSM
- CAA更倾向一些结合不同原理攻击方法
- CAA选择互补的攻击方法加强攻击效果

S_{l_∞}	[(‘MT-LinfAttack’, $\epsilon=8/255$, $t=50$), (‘MT-LinfAttack’, $\epsilon=8/255$, $t=25$), (‘CWLinfAttack’, $\epsilon=8/255$, $t=125$)]
S_{l_2}	[(‘MT-L2Attack’, $\epsilon=0.5$, $t=100$), (‘PGD-L2Attack’, $\epsilon=0.4375$, $t=125$), (‘DDNAttack’, $t=1000$)]
$S_{unrestricted}$	[(‘FogAttack’, $\epsilon = 1$, $t = 1$), (‘FogAttack’, $\epsilon = 1$, $t = 1$), (‘SPSAAttack’, $\epsilon=16/255$, $t=100$)]

- CAA搜索的白盒攻击策略的性能对比实验
 - 数据集
 - CIFAR-10、ImageNet、Bird&Bicycle
 - 实验目的
 - 在开源防御模型上评估CAA搜索到的攻击策略的性能
 - 评价指标
 - 目标模型对生成的对抗样本的准确性: RA (Robust Accuracy)
 - 攻击算法复杂性: Complexity
 - 对比实验设置
 - 对比当前SOTA的攻击方法的效果
 - 对比CAA四种攻击方式的效果
 - Best Attack: 搜索最好的单个攻击(Singe attack)
 - Ens Attack: 搜索最好的攻击组合(Ensemble attack)
 - CAA_{dic}: 直接搜索待攻击模型和数据集
 - CAA_{sub}: 通过攻击对抗训练模型作为替代进行搜索, 并转移到其他模型或任务

- CAA的性能对比实验结果

- 比目前的SOTA: Auto Attack有更好的攻击效果，且运行速度提升6倍，在当前（2020.12.14）白盒算法中取得第一名



- CAA的性能对比实验结果分析
 - Best Attack和Ens Attack没有明显优势
 - CAA_{dic} 和 CAA_{sub} 在增加防御机制的鲁棒模型中的攻击效果要优于目前SOTA的攻击方法，且复杂度更低
 - CAA_{dic} 和 CAA_{sub} 效果相差不大，说明CAA搜索的攻击策略具有可转移性

CIFAR-10 - l_∞ - $\epsilon = 8/255$	AdvTrain	TRADES	AdvPT	MMA	JEM	PCL	Semi-Adv	Complexity
PGD (Madry et al. 2017)	51.95	53.47	57.21	50.04	9.21	8.12	61.83	1000
FAB (Croce and Hein 2019)	49.81	51.70	55.27	42.47	62.71	0.71	60.12	1350
APGD (Croce and Hein 2020)	51.27	53.25	56.76	49.88	9.06	7.96	61.29	1000
AA (Croce and Hein 2020)	49.25	51.28	54.92	41.44	8.15	0.28	59.53	4850
ODI-PGD (Tashiro 2020)	49.37	51.29	54.94	41.75	8.62	0.53	59.61	3200
BestAttack on S_{l_∞}	50.12	52.01	55.23	41.85	9.12	0.84	60.74	900
EnsAttack on S_{l_∞}	49.58	51.51	55.02	41.56	8.33	0.73	60.12	800
CAA_{sub} on S_{l_∞}	49.18	51.19	54.82	40.87	7.47	0.0	59.45	800
CAA_{dic} on S_{l_∞}	49.18	51.10	54.69	40.69	7.28	0.0	59.38	-

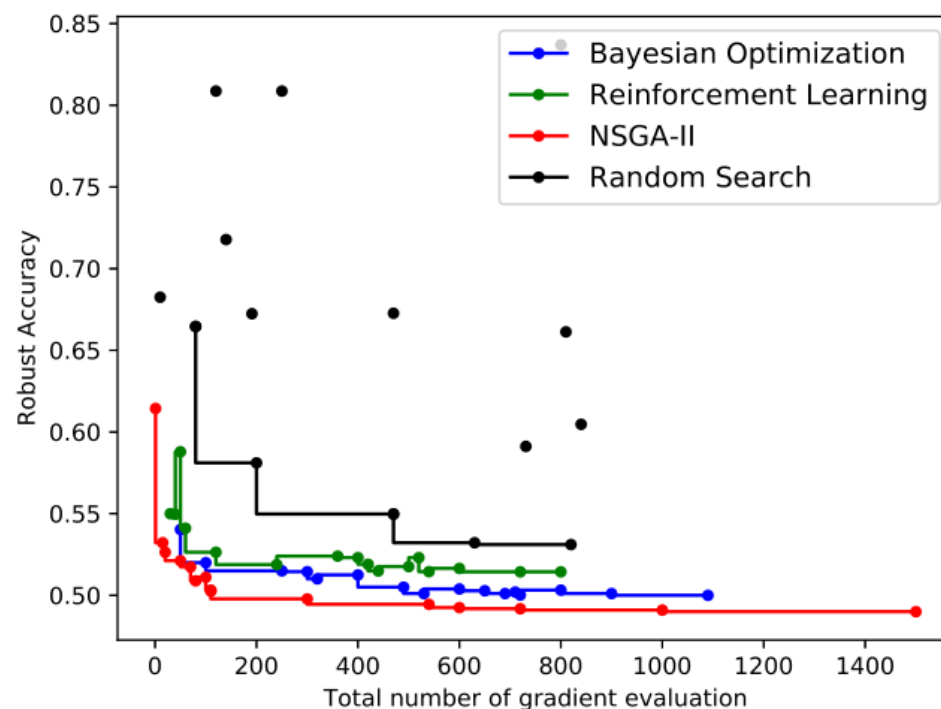
- CAA搜索黑盒攻击策略实验
 - 黑盒攻击：无法获得目标模型的梯度
 - 使用CAA来搜索替代模型上的攻击策略，生成对抗样本来攻击目标模型
 - R→V：使用ResNet50作为替代模型攻击VGG16
 - 结论：CAA也有助于搜索黑盒攻击策略，而**不限于白盒场景**

Models	BestAttack	EnsAttack	CAA _{dic}
R → V	64.75	63.93	63.85
R → I	67.34	66.05	66.21
V → R	67.21	65.23	64.98
V → I	63.42	61.33	60.81
I → R	65.29	64.98	64.38
I → V	59.82	58.54	58.32

- 消融实验1-不同攻击序列长度
- 目的
 - 探求不同攻击序列长度对攻击策略效果的影响
- 对比设置
 - 选择不同的序列长度参数 N ($N = 1, 2, 3, 5, 7$)
- 结论
 - $N = 1$ 时，效果最差
 - 随着 N 的增加，搜索出的攻击策略更好
 - 对比有限制和无限制攻击，攻击序列长度 N 对无限制攻击的性能影响更大
 - 当 N 大于3时，在不受限制的设置中，准确性迅速下降到零左右

- 消融实验2-不同搜索算法
- 目的
 - 衡量时间成本和性能
- 基线
 - 随机搜索100个随机试验策略，选择最佳策略，作为基线
- 对比算法
 - 贝叶斯优化、强化学习、NSGA-II
- 结论
 - 对比NSGA-II，贝叶斯优化和强化学习算法，更费时，性能略差

Search Methods	Performance	Search time
Random Search-100	52.09	8 Hours
Reinforcement Learning	51.44	5 GPU/d
Bayesian Optimization	50.02	5 GPU/d
NSGA-II	49.18	3 GPU/d





优劣分析

- 优势

- 通过自动化搜索攻击序列组合和超参数设置，减少人工干预
- 搜索的最佳攻击序列组合提高了对抗攻击性能
- 相比于之前的SOTA方法减少计算成本

- 劣势

- 搜索空间的设置困难（攻击候选池、参数搜索范围）
- 目前仅限应用于图片分类任务

- 对抗样本攻击自动化搜索在智能系统测试中的应用
 - 黄、赌、毒等不良内容的安全检测系统
 - 曝光、模糊、低画质等极端分类场景

[(ContrastAttack, $\epsilon = 1$, $t = 1$), (ShotNoiseAttack, $\epsilon = 1$, $t = 1$), (SPSAAttack, $\epsilon = 16/255$, $t = 100$)]



[(ContrastAttack, $\epsilon = 1$, $t = 1$), (FogAttack, $\epsilon = 1$, $t = 1$), (SPSAAttack, $\epsilon = 16/255$, $t = 100$)]



- [1] Mao X , Chen Y , Wang S , et al. Composite Adversarial Attacks[J]. 2020.
- [2] Elshaw R , Maher M , Sakr S . Automated Machine Learning: State-of-The-Art and Open Challenges[J]. 2019.
- [3] Xu H , Ma Y , Liu H C , et al. Adversarial Attacks and Defenses in Images, Graphs and Text: A Review[J]. International Journal of Automation and Computing, 2020, 17(2):151–178.
- [4] Akhtar N , Mian A . Threat of Adversarial Attacks on Deep Learning in Computer Vision: A Survey[J]. IEEE Access, 2018, 6:14410–14430.

知人者智，自知者明。
胜人者有力，自胜者
强。知足者富。强行
者有志。不失其所者
久。死而不亡者，寿。

谢谢！

