



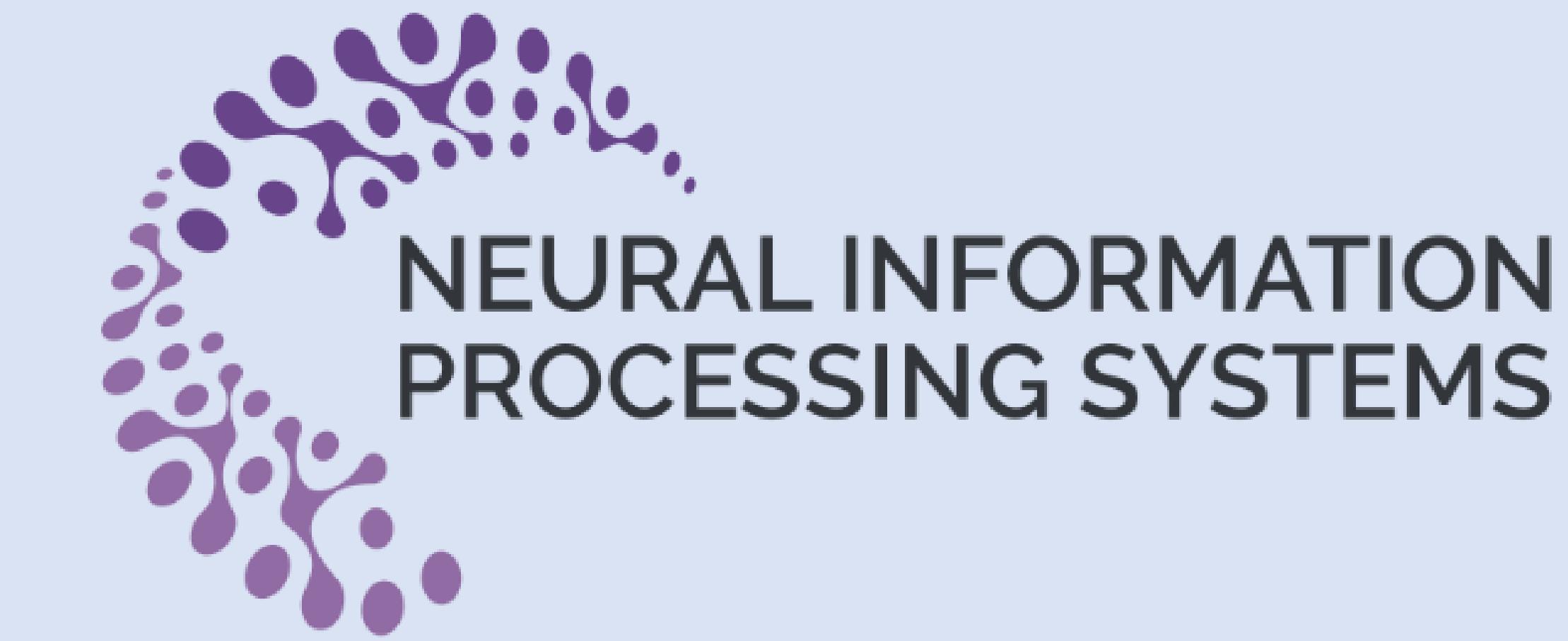
Finding Optimal Tangent Points for Reducing Distortions of Hard-label Attacks

Chen Ma¹, Xiangyu Guo², Li Chen¹, Jun-Hai Yong¹, Yisen Wang³

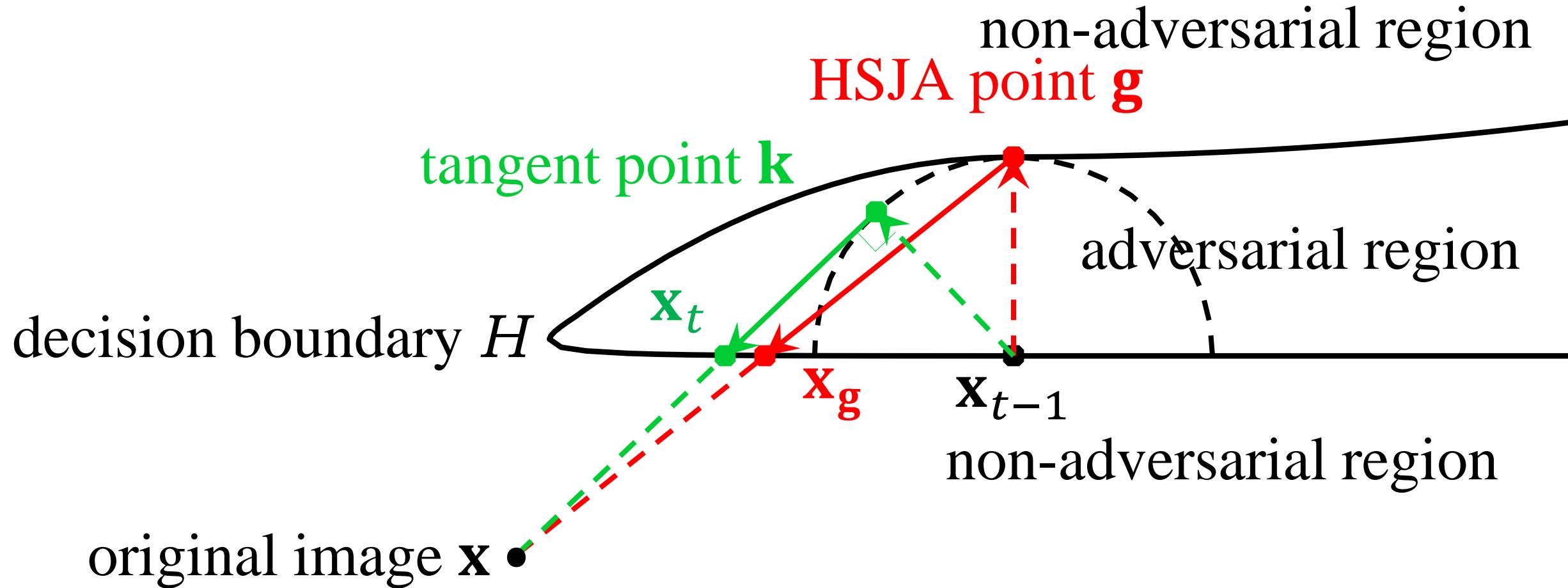
¹ School of Software, BNRIst, Tsinghua University, Beijing, China

² Department of Computer Science and Engineering, University at Buffalo, Buffalo NY, USA

³ Key Lab. of Machine Perception, School of EECS, Peking University, Beijing, China



Introduction



The hard-label attack is a challenging adversarial attack, in which only the top-1 predicted label is available. Existing hard-label attacks do not thoroughly investigate the geometric properties of the decision boundary to accelerate the attack. Let us take HSJA (HopSkipJumpAttack) for example. As shown in the above figure, x_{t-1} is the current adversarial example mapped onto the decision boundary at the $(t-1)$ -th iteration of the attack. HSJA updates x_{t-1} along the gradient direction to reach g and then maps it to H at x_g along the line through x and g . However, the optimal update is not g . It is easy to observe that moving along the tangent line can reach the nearest location of the decision boundary to the original image x , thereby producing the adversarial example with the minimum distortion. In real attack scenarios, the image data reside in a high-dimensional space, and the semicircle becomes a hemisphere B . In n -dimensional space where $n \geq 3$, there are infinitely many tangent lines from x to B which produce infinitely many tangent points on B . Still, we will show that **exactly one tangent point** can lead to the minimum distortion when mapping it onto H along the tangent line, which is defined by the following theorem.

Definition of Optimal Tangent Points

Theorem 1: Let u be the unit normal vector of H , k be any point on the surface of B , then the distance $\|x - x_t\|_2$ is the shortest if k is the optimal solution of the following constrained optimization problem:

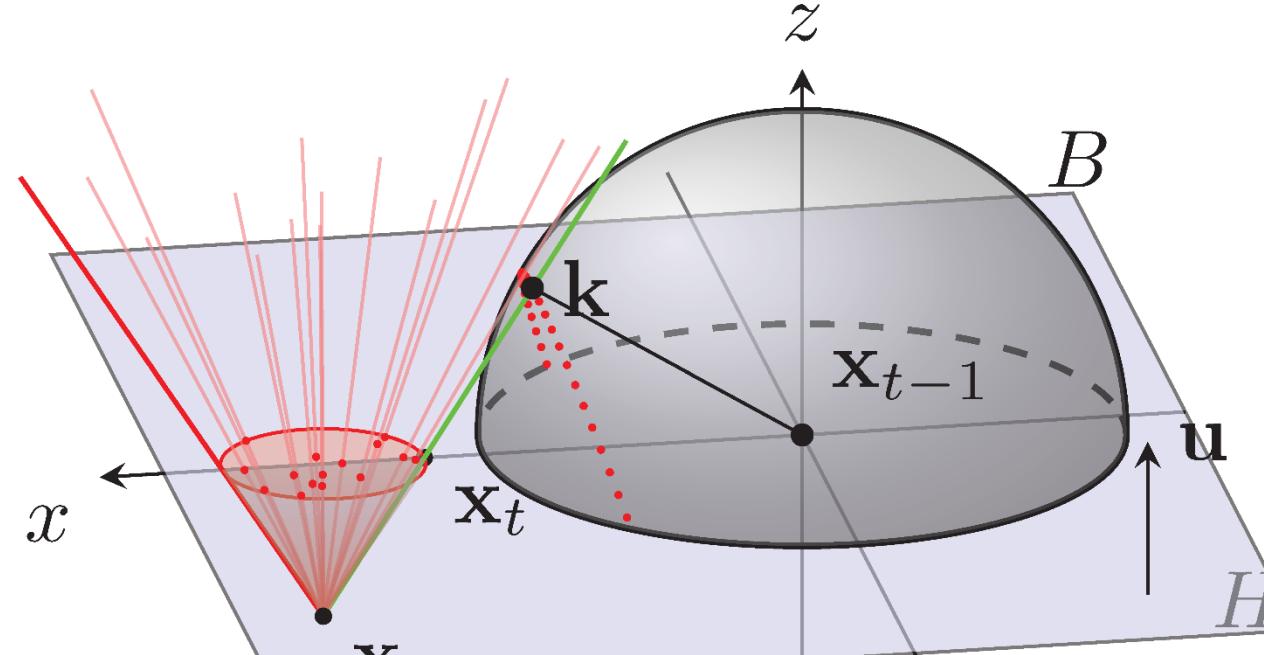
$\underset{k}{\operatorname{argmax}} \langle k - x_{t-1}, u \rangle$, —— maximizes the projection of $k - x_{t-1}$ onto u .

s.t. $\langle k - x_{t-1}, x - k \rangle = 0$, —— ensures that k is a tangent point.

$\|k - x_{t-1}\| = 0$, —— indicates k is on the surface of B .

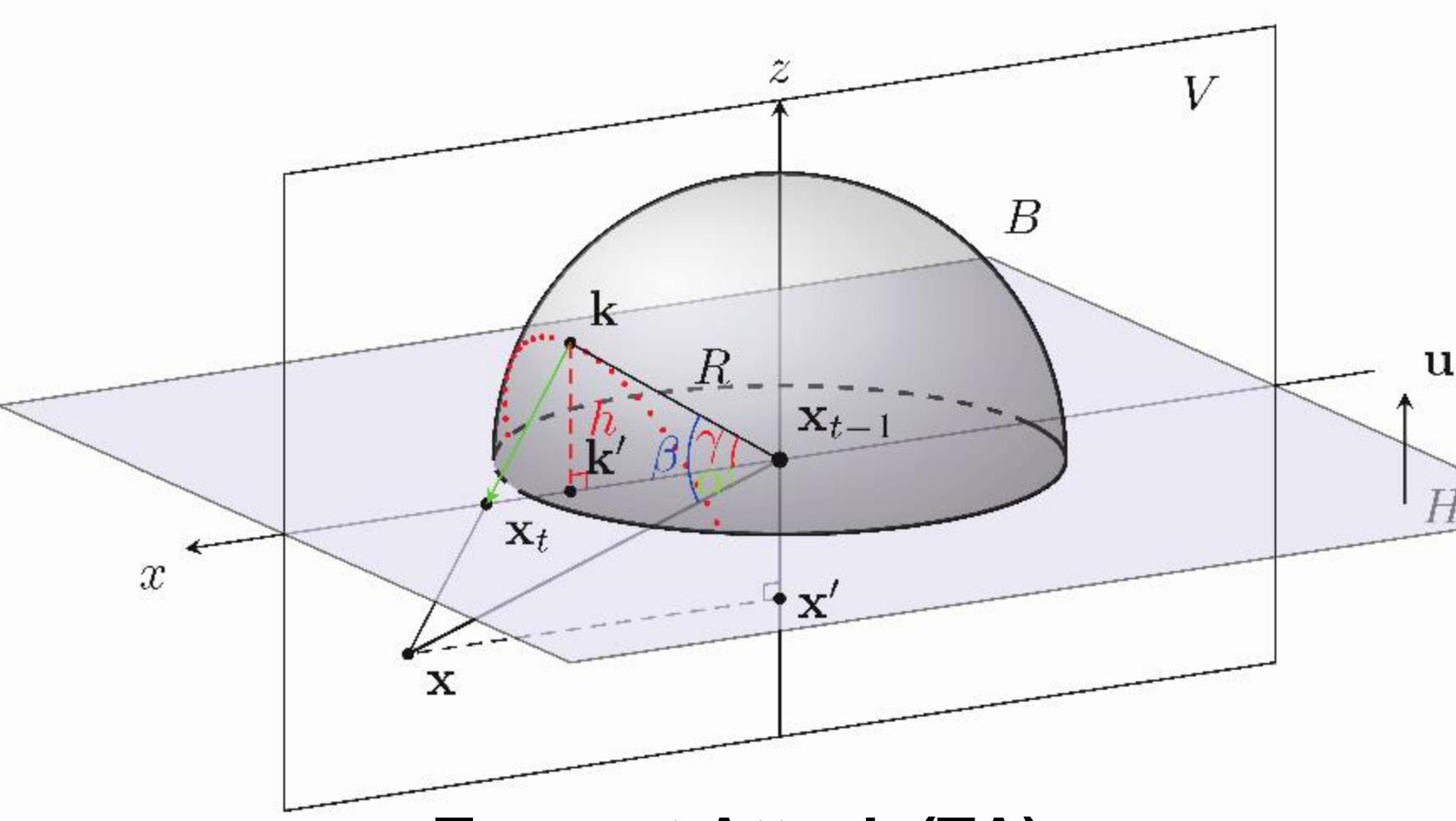
$\langle k - x_{t-1}, u \rangle > 0$, —— states that k cannot appear on the same side of H as x , which is always satisfied in our algorithm.

Geometrical Explanation of Theorem 1



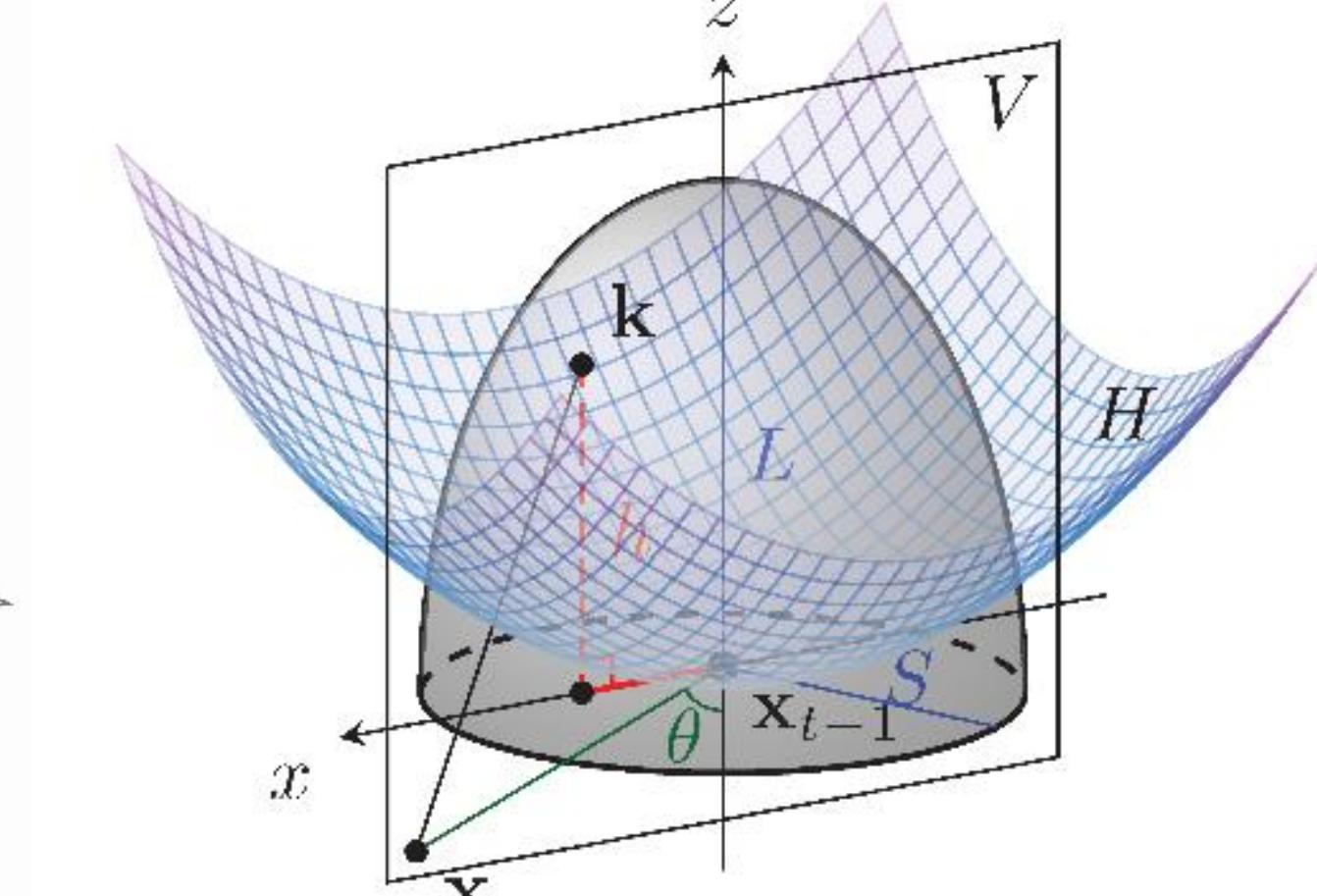
The left figure shows that all points on H that are closer to x than x_t are within a red disk, which is the intersection of H and the red cone whose vertex is x . Clearly k is the only intersection point of the cone and the hemisphere B . Thus, of all the lines intersecting B , only the tangent line leads to the shortest distance from x to H .

Computation of Optimal Tangent Points



Tangent Attack (TA)

$$\begin{aligned} \beta &= \alpha + \gamma \\ \sin \alpha &= -\frac{\langle x, u \rangle}{\|x\| \cdot \|u\|} \\ \cos \alpha &= \sqrt{\|x\|^2 \cdot \|u\|^2 - \langle x, u \rangle^2} \\ \sin \beta &= \sqrt{1 - \cos^2 \beta} \\ \cos \gamma &= \cos \beta \cos \alpha + \sin \beta \sin \alpha \\ \sin \gamma &= \sin \beta \cos \alpha - \cos \beta \sin \alpha \\ h &= R \cdot \sin \gamma = R \cdot (\sin \beta \cos \alpha - \cos \beta \sin \alpha) \\ x' &= \langle x, u \rangle \cdot u / \|u\|^2 \\ k &= k' + h \cdot u \\ &= R \cdot \cos \gamma \cdot \frac{x - \langle x, u \rangle u / \|u\|^2}{\|x - \langle x, u \rangle u / \|u\|^2\|} + h \cdot u \end{aligned}$$

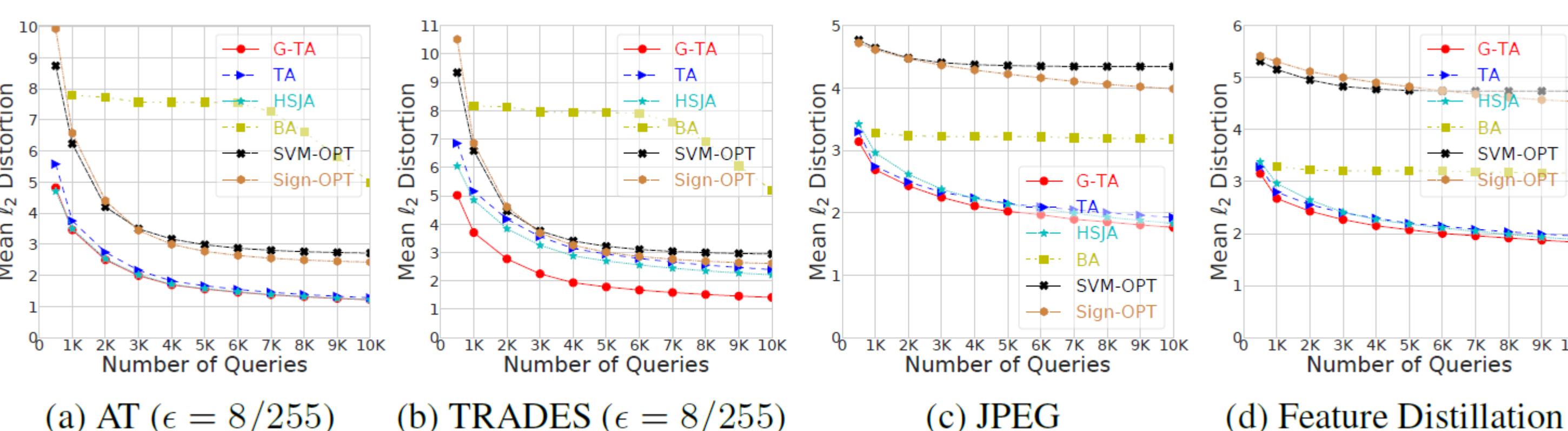


Generalized Tangent Attack (G-TA)

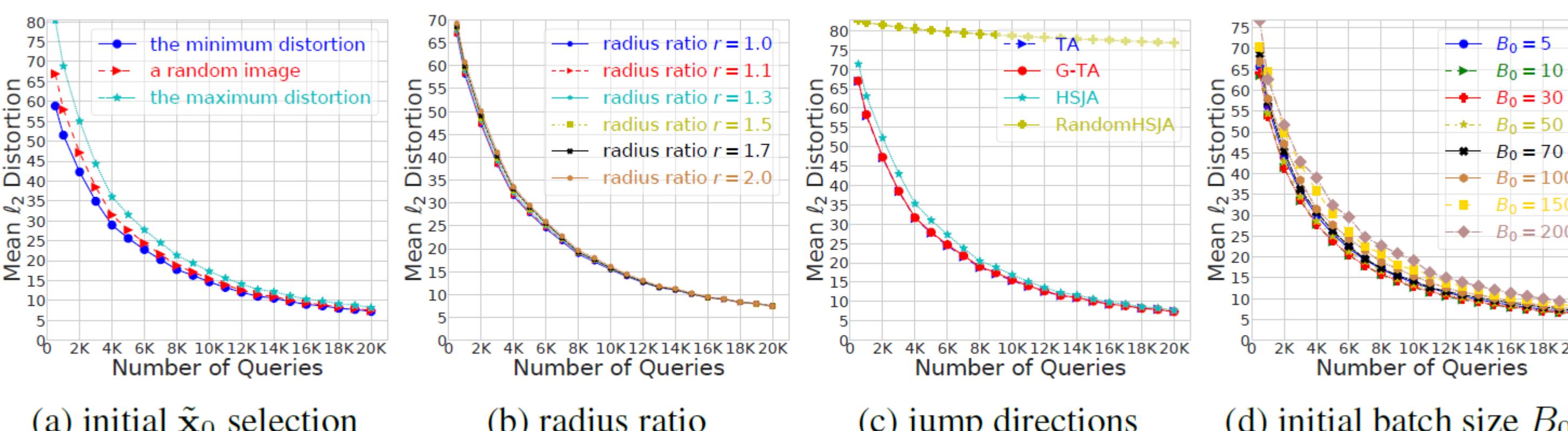
$$\begin{aligned} \theta &= \arccos \left(\frac{\langle x, u \rangle}{\|x\| \cdot \|u\|} \right) \\ (x_0, z_0) &= (\|x\| \cdot \sin \theta, -\|x\| \cdot \cos \theta) \\ x_k &= \frac{s^2 \left(L^2 z_0 + L^2 x_0 \sqrt{-L^2 s^2 + L^2 x_0^2 + s^2 z_0^2} \right)}{L^2 x_0 + s^2 z_0^2} \\ z_k &= \frac{L^2 s^2 z_0 + L^2 x_0 \sqrt{-L^2 s^2 + L^2 x_0^2 + s^2 z_0^2}}{L^2 x_0^2 + s^2 z_0^2} \end{aligned}$$

$$k = |x_k| \cdot \frac{x - \langle x, u \rangle u / \|u\|^2}{\|x - \langle x, u \rangle u / \|u\|^2\|} + z_k \cdot u$$

Experimental Results of Attacks against Defense Models



Comprehensive Understanding of Tangent Attack



Experimental Results

Table 1: Mean ℓ_2 distortions of different query budgets on the ImageNet dataset, where $r = 1.1$.

Target Model	Method	Targeted Attack					Untargeted Attack					
		@300	@1K	@2K	@5K	@8K	@10K	@300	@1K	@2K	@5K	
Inception-v3	BA [4]	111.798	108.044	106.283	102.715	86.931	78.326	-	107.558	102.309	95.776	78.668
	Sign-OPT [11]	103.939	87.706	71.291	46.744	34.640	29.414	121.085	79.158	43.642	16.625	10.557
	SVM-OPT [11]	101.630	82.950	67.965	46.275	35.694	31.106	121.135	66.027	36.763	15.736	10.501
	HSJA [7]	111.562	95.295	82.111	52.544	37.395	30.425	103.605	57.295	37.185	15.484	9.989
	TA	103.781	80.327	66.708	42.121	30.846	25.566	94.752	52.523	35.229	15.040	9.748
Inception-v4	G-TA	103.724	81.089	67.168	42.434	31.011	25.587	94.668	52.037	34.540	14.643	9.540
	BA [4]	110.343	106.616	104.586	100.321	84.058	75.507	-	116.075	111.474	104.451	86.572
	Sign-OPT [11]	101.620	85.731	69.719	46.416	34.957	30.004	132.991	86.431	48.292	18.678	11.567
	SVM-OPT [11]	99.856	81.342	66.982	45.667	35.477	31.152	132.227	72.920	41.095	17.611	11.418
	HSJA [7]	109.670	93.916	80.937	52.358	37.773	30.958	110.727	63.731	42.290	17.936	11.367
SENet-154	TA	101.666	78.683	65.304	41.629	30.993	25.958	101.207	58.616	40.314	17.639	11.304
	G-TA	101.495	79.210	65.888	42.002	30.965	25.847	101.173	58.225	39.788	17.265	11.008
	BA [4]	81.090	77.723	76.122	71.967	55.953	47.652	-	75.998	71.671	66.983	53.917
	Sign-OPT [11]	75.722	62.876	49.191	30.155	21.333	17.672	70.035	47.705	27.314	10.890	6.245
	SVM-OPT [11]	74.658	58.677	46.827	30.264	22.461	19.186	69.854	40.291	23.692	10.494	6.666
ResNet-101	HSJA [7]	77.035	63.488	51.802	30.138	19.680	16.261	71.248	38.035	24.895	10.218	5.855
	TA	70.739	55.256	43.694	24.961	16.756	13.876	65.589	35.689	24.037	10.039	5.774
	G-TA	70.591	55.224	44.047	25.041	16.854	14.047	65.871	35.768	23.954	9.959	5.733
	BA [4]	81.565	77.903	76.366	72.392	58.746	51.679	-	64.007	60.389	56.544	44.175
	Sign-OPT [11]	76.732	63.939	51.231	32.439	23.160	7.050	56.244	38.282	21.985	10.048	7.050
ResNet-101	SVM-OPT [11]	77.031	61.417	49.842	32.806	24.553	20.964	55.894	32.638	19.409	9.830	7.185
	HSJA [7]	76.121	63.091	52.301	31.018	20.472	16.911	56.264	27.443	17.717	7.649	4.723
	TA	72.434	57.969	47.142	27.699	18.788	15.414	53.197	26.777	17.651	7.730	4.822
	G-TA	72.459	58.320	47.297	27.905	19.045	15.633	53.058	26.631	17.384	7.602	4.720
	BA [4]	81.565	77.903	76.366	72.392	58.746	51.679	-	64.007	60.389	56.544	44.175

Table 2: Mean ℓ_2 distortions with different query budgets on the CIFAR-10 dataset, where $r = 1.5$.