# Towards Scalable Coverage-Based Testing of Autonomous Vehicles - Supplementary Materials

**Anonymous Author(s)**
Affiliation
Address
`email`

## 1 Additional Experiments

**Ablation on GP Kernels:** While GUARD requires minimal tuning of GP kernel parameters due to online kernel learning, the *type* of kernel is still a design choice. In GUARD we employ Matern32 kernels to model each parameter and take the product of individual kernels to obtain the final kernel. In Table 1 we conduct an ablation that also considers 1) RBF kernels, 2) Using sum instead of product, and 3) no composition, using one RBF or Matern32 kernel for all parameters. We first find that without modeling each parameter separately, the results are poor. Also, using a sum operation to aggregate kernels is ineffective as well. The best results are obtained from a product of RBF or Matern32 kernels. We choose product of Matern32 as the lower false positive rate is critical for ensuring safety. Furthermore, this choice also exhibits much lower variance in the metrics, making it more reliable for testing.

**BO Acquisition Functions:** As mentioned in Section 2 of the main manuscript, many adversarial example generation methods use Bayesian Optimization(BO) [1] and employ acquisition functions that minimize some adversarial objective. Since BO also leverages GPs and active learning similar to GUARD, we can adopt these acquisition functions. In particular we use lower confidence bound (LCB), expected improvement (EI), and probability of improvement (PI). For EI and PI we adjust them to decrease the objective (expected decrease and probability of decrease). The results in Table 2 show that these alternatives can achieve strong results but is not as effective as using the Straddle acquisition function. This is because cost minimization oversample very negative points and Straddle oversamples boundary points, but the former is more useful in partitioning pass / fail regions. The performance is however very close, and we hypothesize this is because the parameter space is dominated by positive (passes). With very small negative regions, severe negatives are close to boundary points.

**Runtime Comparison** A major implementation detail regarding our baselines is the discretization resolution. For grid search and $t$-way testing, there is a direct tradeoff between coverage and balanced accuracy, error recall, false positive rate. Thus, we select a resolution which can cover most of the parameter space. For HiddenGems [2], the resolution does not affect coverage since it uses a levelset estimation [3] algorithm to determine which bins are covered. The limitation for our choice of resolution of 6 is due to how the computation budget increases exponentially with resolution. This is because the method requires running the GP on all discretized bins at every iteration. And while running the GP a single time is neglible compared to running simulation, the exponential scaling makes the GP queries a



Figure 1: Runtime

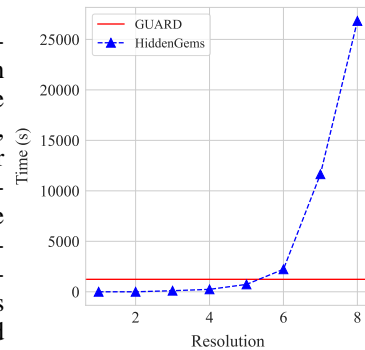bottleneck at higher resolutions. As we show in Figure 1, at this resolution the runtime roughly equals GUARD.

## 2 Algorithm Details

We provide a more detailed overview of the algorithm used for GUARD in Algorithm 1. The $\mathrm{argmax}$ operations are done via gradient based optimization as described in the main manuscript.

| Kernel Type | Composition | Coverage(%) | Bal. Acc(%) | Pos Acc(%) | Neg Acc(%) | Err. Recall(%) | FPR(%) |
|---|---|---|---|---|---|---|---|
| RBF | - | $96.1 \pm 0.8$ | $80.2 \pm 5.7$ | $98.1 \pm 0.5$ | $62.4 \pm 11.2$ | $56.0 \pm 10.6$ | $10.6 \pm 3.3$ |
| Matern32 | - | $87.3 \pm 1.5$ | $78.5 \pm 4.2$ | $\mathbf{99.3 \pm 0.2}$ | $57.7 \pm 8.4$ | $36.8 \pm 7.7$ | $6.15 \pm 1.1$ |
| RBF | Sum | $88.0 \pm 6.7$ | $58.6 \pm 3.2$ | $89.1 \pm 11.4$ | $28.2 \pm 11.5$ | $21.4 \pm 12.0$ | $20.2 \pm 5.7$ |
| Matern32 | Sum | $90.4 \pm 3.3$ | $58.4 \pm 3.3$ | $88.0 \pm 10.7$ | $28.7 \pm 13.7$ | $23.6 \pm 12.7$ | $21.0 \pm 3.1$ |
| RBF | Product | $96.8 \pm 0.7$ | $83.3 \pm 7.0$ | $97.9 \pm 0.6$ | $68.7 \pm 13.9$ | $\mathbf{62.6 \pm 13.1}$ | $7.83 \pm 4.0$ |
| Matern32 | Product | $94.3 \pm 0.9$ | $\mathbf{84.2 \pm 5.3}$ | $99.2 \pm 0.2$ | $70.3 \pm 10.6$ | $58.9 \pm 5.3$ | $\mathbf{5.89 \pm 3.1}$ |

Table 1: Ablation of different kernel choices for GUARD

| Aquisition Fn. | Coverage(%) | Bal. Acc(%) | Pos Acc(%) | Neg Acc(%) | Err. Recall(%) | FPR(%) |
|---|---|---|---|---|---|---|
| PI | $90.7 \pm 1.6$ | $82.2 \pm 3.9$ | $99.4 \pm 0.1$ | $64.9 \pm 7.8$ | $49.0 \pm 7.5$ | $8.23 \pm 2.4$ |
| EI | $91.8 \pm 1.3$ | $82.6 \pm 2.3$ | $99.4 \pm 0.2$ | $65.7 \pm 0.5$ | $50.7 \pm 5.2$ | $6.99 \pm 1.4$ |
| LCB | $91.5 \pm 1.5$ | $83.4 \pm 1.8$ | $\mathbf{99.6 \pm 0.1}$ | $67.3 \pm 3.6$ | $52.4 \pm 4.7$ | $5.99 \pm 1.4$ |
| Straddle | $\mathbf{94.3 \pm 0.9}$ | $\mathbf{84.2 \pm 5.3}$ | $99.2 \pm 0.2$ | $\mathbf{70.3 \pm 10.6}$ | $\mathbf{58.9 \pm 5.3}$ | $\mathbf{5.89 \pm 3.1}$ |

Table 2: Comparison of cost minimization acquisition functions versus Straddle acquistion function.

---

**Algorithm 1:** Detailed Algorithm for GUARD

**Input:** Test function $f^*$, test Threshold $\gamma$, budget $N$, initial batch size $M$, initial kernel parameters $\boldsymbol{\kappa}_0$, kernel update frequency $K$, Gaussian CDF $\Phi(\cdot)$

$G \leftarrow GP(\boldsymbol{\kappa}_0)$      ▷ initialize GP

$\mu(\cdot), \sigma(\cdot) \leftarrow$ mean of $G$, standard deviation of $G$

**for** $t$ *in 1, ..., M* **do**
    $\boldsymbol{\theta}_t \leftarrow \arg\max_{\boldsymbol{\theta}} \sigma(\boldsymbol{\theta})$      ▷ Initial batch for exploration
    $G \leftarrow \text{UPDATEGPSAMPLES}(G, \boldsymbol{\theta}_t, f^*(\boldsymbol{\theta}_t))$
**end**

$\boldsymbol{\kappa}_M \leftarrow \arg\max_{\boldsymbol{\kappa}} P(f(\boldsymbol{\theta}_1) \dots f(\boldsymbol{\theta}_t) \mid \boldsymbol{\theta}_1 \dots \boldsymbol{\theta}_t; \boldsymbol{\kappa})$      ▷ Probability under GP prior

$G \leftarrow \text{UPDATEGPKERNEL}(G, \boldsymbol{\kappa}_M)$

**for** $t$ *in M+1, ..., N* **do**
    $\boldsymbol{\theta}_t \leftarrow \arg\max_{\boldsymbol{\theta}} \beta\sigma(\boldsymbol{\theta}) + |\mu(\boldsymbol{\theta}) - \gamma|$      ▷ Explore or be close to boundary
    $G \leftarrow \text{UPDATEGP}(G, \boldsymbol{\theta}_t, f^*(\boldsymbol{\theta}_t))$
    **if** $t \equiv 0 \mod K$ **then**
        $\boldsymbol{\kappa}_t \leftarrow \arg\max_{\boldsymbol{\kappa}} P(f(\boldsymbol{\theta}_1) \dots f(\boldsymbol{\theta}_t) \mid \boldsymbol{\theta}_1 \dots \boldsymbol{\theta}_t; \boldsymbol{\kappa})$ $G \leftarrow \text{UPDATEGPKERNEL}(G, \boldsymbol{\kappa}_t)$
    **end**
**end**

**function** PASSFAILESTIMATION$(\boldsymbol{\theta}, \alpha)$ **is**
    $p^+ \leftarrow \Phi\left(\frac{\mu(\boldsymbol{\theta}) - \gamma}{\sigma(\boldsymbol{\theta})}\right)$      ▷ Probability value at $\boldsymbol{\theta}$ is greater than threshold

    $p^- \leftarrow \Phi\left(\frac{\gamma - \mu(\boldsymbol{\theta})}{\sigma(\boldsymbol{\theta})}\right)$      ▷ Probability value at $\boldsymbol{\theta}$ is less than threshold

    **if** $p^+ \geq \alpha$ **then**
        **return** 1      ▷ Pass
    **else if** $p^- \geq \alpha$ **then**
        **return** $-1$      ▷ Fail
    **else**
        **return** 0      ▷ Uncertain
**end**

**return** PASSFAILESTIMATION

---

## 3 Test Scenarios

Details on our parameterized scenarios are provided in Table 3 where we give a high level description of each scenario, associated parameters configuring the traffic agents or road, and parameter bounds.

| Scenario | Parameter | Bounds |
|---|---|---|
| AV follows lane, actor cuts in from neighboring lane | AV initial speed<br>Actor initial speed relative to AV<br>Cut-in lane change duration<br>Time to collision at initialization assuming constant speed<br>Road curvature | [20, 40] m/s<br>[-10, 10] m/s<br>[1, 10] s<br>[1, 6] s<br>[-0.002, 0.002] $m^{-1}$ |
| AV follows lane, multiple actors merge in from on-ramp | AV time-to-arrival to merge point<br>AV initial speed relative to speed limit<br>Actors speed relative to AV<br>Time interval between actors<br>Actors time-to-arrival to merge point | [0.5, 5.0] s<br>[-10.0, 10.0] m/s<br>[-10.0, 10.0] m/s<br>[1.0, 4.0] m/s<br>[0.5, 5.0] s |
| AV follows lane, lead actor brakes | Time-to-collision at initialization assuming constant speed<br>Actor initial distance relative to AV<br>AV initial speed relative to speed limit<br>Actor target speed relative to AV initial speed<br>Road curvature | [1.0, 10.0] s<br>[20.0, 200.0] m<br>[-10.0, 0.0] m/s<br>[-30.0, 10.0] m/s<br>[-0.003, 0.003] $m^{-1}$ |
| AV follows lane, actor overtakes from neighboring lane | Actor initial follow distance<br>Actor headway for lane change<br>Actor overtake acceleration<br>Actor lane change duration<br>Road curvature | [1.0, 30.0] m/s<br>[-10.0, 10.0] m<br>[0.1, 5.0] $m/s^2$<br>[0.5, 3.0] s<br>[-0.002, 0.002] $m^{-1}$ |
| AV follows lane, actor cuts in from shoulder | Actor initial speed relative to AV<br>Actor initial distance from lane boundary<br>Time-to-collision at initialization assuming constant speed<br>Actor cut-in duration<br>Road curvature | [-29.0, -15.0] m/s<br>[0.1, 0.3] m<br>[1.0, 6.0] s<br>[1.0, 8.0] s<br>[-0.002, 0.002] $m^{-1}$ |
| AV merges from on-ramp, with an actor on its target lane | AV time-to-arrival to merge point<br>Actor time-to-arrival to merge point<br>AV initial speed relative to speed limit<br>Actor initial speed relative to AV initial speed<br>Lane encroachment threshold for actor to react to AV | [0.5, 6.0] s<br>[0.5, 6.0] s<br>[-10.0, 10.0] m/s<br>[-10.0, 10.0] m/s<br>[-5.0, 5.0] m |
| AV merges from parallel on-ramp, with an actor on its target lane | AV time-to-arrival to merge point<br>Actor time-to-arrival to merge point<br>AV initial speed relative to speed limit<br>Actor initial speed relative to AV initial speed<br>Lane encroachment threshold for actor to react to AV | [0.5, 6.0] s<br>[0.5, 6.0] s<br>[-10.0, 10.0] m/s<br>[-10.0, 10.0] m/s<br>[-5.0, 5.0] m |
| AV merges from on-ramp, with multiple actors on its target lane | AV time-to-arrival to merge point<br>Actors time-to-arrival to merge point<br>AV initial speed relative to speed limit<br>Actors initial speed relative to AV initial speed<br>Time interval between actors | [0.5, 6.0] s<br>[-8.0, 4.0] s<br>[-10.0, 10.0] m/s<br>[-10.0, 10.0] m/s<br>[1.0, 6.0] s |
| AV changes lane, with an actor on its target lane | Road curvature<br>Actor initial speed relative to AV initial speed<br>AV initial speed relative to speed limit<br>Time-to-collision at initialization assuming constant speed<br>Lane encroachment threshold for actor to react to AV | [-0.002, 0.002] $m^{-1}$<br>[-10.0, 10.0] m/s<br>[-10.0, 10.0] m/s<br>[1.0, 15.0] s<br>[-2.0, 2.0] m |
| AV changes lane, with multiple actors merging from on-ramp | AV time-to-arrival to merge point<br>Actor time-to-arrival to merge point<br>AV initial speed relative to speed limit<br>Actor initial speed relative to AV initial speed<br>Time interval between actors | [0.5, 6.0] s<br>[0.5, 6.0] s<br>[-10.0, 10.0] m/s<br>[-10.0, 10.0] m/s<br>[1.0, 6.0] s |

Table 3: Test scenario details

## References

[1] A. Agnihotri and N. Batra. Exploring bayesian optimization. *Distill*, 2020. doi:10.23915/distill.00026. https://distill.pub/2020/bayesian-optimization.

[2] A. Petrov, C. Fang, K. M. Pham, Y. H. Eng, J. G. M. Fu, and S. D. Pendleton. Hiddengems: Efficient safety boundary detection with active learning. *arXiv preprint arXiv:2210.13956*, 2022.

[3] A. Gotovos. Active learning for level set estimation. Master's thesis, Eidgenössische Technische Hochschule Zürich, Department of Computer Science,, 2013.