
Beyond Dice: Risk-Normalized and Hazard-Aware Evaluation of Medical Segmentation for Image-Guided Robotics

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Safety in embodied systems depends on where perception fails, not only on how
2 often it fails. We introduce two bounded evaluation metrics for segmentation
3 that make spatial risk explicit using an anatomy-derived hazard field built from
4 distance to protected structures. The *Safety Impact Score* (SIS) measures the share
5 of total hazard mass that is misclassified, with a tunable trade-off between false
6 negatives and false positives. The *Safety Tail Risk* (STAR) summarizes the worst
7 fraction of error hazards using a conditional value-at-risk operator. To isolate
8 metric behavior from model quality, we design a model-free matched-Dice stress
9 test that relocates equal numbers of boundary errors toward or away from hazard
10 while keeping Dice unchanged. We run this protocol on three public Medical
11 Segmentation Decathlon tasks (Hepatic Vessel, Liver, Pancreas; five cases each).
12 Across datasets, STAR shows large positive deltas for the risky variant (combined
13 mean $\Delta\text{STAR} = 0.431 \pm 0.124$, paired Wilcoxon $p = 3.24 \times 10^{-4}$), and SIS is
14 also positive (combined mean $\Delta\text{SIS} = 0.0818 \pm 0.0771$, $p = 3.05 \times 10^{-5}$). Effects
15 are strongest when the hazard corresponds to vessels in Hepatic Vessel (mean
16 $\Delta\text{STAR} = 0.614 \pm 0.307$, mean $\Delta\text{SIS} = 0.224 \pm 0.182$). A proximity-weighted
17 overlap baseline (hazard-weighted Dice) moves little or in the opposite direction.
18 Results persist with an exponential hazard kernel, indicating robustness. These
19 findings demonstrate that risk-normalized and tail-aware evaluation captures safety-
20 relevant differences that overlap metrics miss, using only public data and a simple
21 perturbation protocol.

22 1 Background

23 1.1 Why average segmentation accuracy is not safety

24 Perception modules increasingly guide actions in image-guided and robotic procedures, where the
25 severity of an error depends on *where* it occurs. A small miss that pushes a boundary toward a
26 critical structure can be far more consequential than a larger miss in benign tissue. Conventional
27 metrics such as Dice, Jaccard, and boundary distances treat all voxels or boundary points roughly
28 uniformly and typically report means or percentiles aggregated over space and cases. Multiple studies
29 show that these scores can fail to predict downstream harm or workload. In radiotherapy, geometric
30 scores correlate weakly with dose deviations and clinical acceptability, which motivates multi-aspect
31 evaluation that includes treatment impact and expert ratings rather than geometry alone Poel et al.
32 (2021, 2025); Maier-Hein et al. (2024); Reinke et al. (2023). Work on edit-effort and clinical usability
33 points the same way: Added Path Length, surface-tolerant measures such as Normalized Surface
34 Dice, and the Mendability Index track how hard clinicians must work to fix outputs, not just how

35 similar masks look He et al. (2024); Various (2024); Ma et al. (2024); He et al. (2023); Zhang et al.
36 (2024). Average overlap is therefore a poor proxy for *risk*.

37 1.2 What has already been done to inject safety into metrics

38 **Proximity- or sensitivity-weighted overlap in radiotherapy.** Several groups weight overlap by
39 anatomical importance. OAR-DSC reweights Dice by proximity to organs at risk and by radiosensi-
40 tivity McCullum et al. (2024). Weighted Dice for brachytherapy similarly uses distance to high-risk
41 volumes and reports stronger links to dose endpoints Ni et al. (2025). These are important steps
42 toward task-aware evaluation.

43 **Outcome-aware validation frameworks.** Comprehensive frameworks evaluate contouring by
44 geometry, efficiency, plan quality, and expert acceptability, which exposes outliers missed by average
45 metrics Poel et al. (2025). Guidance now recommends reporting complementary measures and
46 considering task risk explicitly Maier-Hein et al. (2024); Reinke et al. (2023).

47 **Tail-sensitive boundary metrics.** Hausdorff distance and HD95 emphasize extreme boundary
48 errors Huttenlocher et al. (1993); Karimi & Salcudean (2019). They are useful for spotting worst-case
49 deviations, yet agnostic to which anatomy those deviations threaten.

50 **Safety-aware perception outside medicine.** In autonomous driving, safety metrics for semantic
51 segmentation penalize clustered errors and errors in safety-critical regions instead of averaging over
52 all pixels Cheng et al. (2021). Task-aware risk estimators quantify how perception failures propagate
53 into risky plans and highlight low-probability high-impact events Antonante et al. (2023).

54 **Risk modeling around vital anatomy for robotic procedures.** Surgical robotics often frames
55 safety as distance-to-hazard with safety margins that expand where uncertainty is higher. Inner ear
56 and bone-milling work formalize probabilistic error budgets and turn them into spatial keep-out
57 regions around critical structures Dillon et al. (2016); Siebold et al. (2017). This connects naturally
58 to evaluating perception in terms of *clearance* to protected anatomy.

59 1.3 Where the gaps remain

- 60 • **Granular spatial risk in the metric itself.** Existing weighted Dice variants attach scalar
61 weights to overlap, but they do not evaluate errors through a continuous hazard field that
62 reflects millimeter-scale clearance to protected structures. HD95 is tail-sensitive but not
63 hazard-aware.
- 64 • **Risk-mass normalization for cross-case comparability.** Most proposals reweight a
65 score rather than compute the *share of available risk* that the model gets wrong. Without
66 normalizing by total hazard mass, the same numeric change can mean very different things
67 across patients or scanners.
- 68 • **False negative versus false positive asymmetry.** Few evaluation metrics expose explicit
69 cost trade-offs between risk-weighted FN and FP, even though misses near critical anatomy
70 are often costlier than conservative over-segmentations.
- 71 • **Tail risk over hazard, not just over distance.** There is no standard evaluation that
72 aggregates the *worst fraction of errors by hazard* in the spirit of conditional value at risk.
73 Percentile Hausdorff is a boundary percentile, not a hazard-weighted tail of clinically
74 dangerous errors.

75 1.4 How our approach is different

76 We introduce a small, principled family of evaluation metrics that make safety the first-class quantity.

77 **Spatial hazard map from anatomy.** We build a continuous hazard field $w(x) \in [0, 1]$ from
78 distance to critical structures, with polynomial or exponential decay. The field can combine multiple
79 organs at risk via conservative max aggregation or additive accumulation.

80 **Risk-normalized error shares.** Let $y(x) \in \{0, 1\}$ be ground truth for the target, $\hat{y}(x) \in \{0, 1\}$ a
 81 prediction, and $w(x)$ the hazard map. Define

$$\text{R-FN} = \frac{\sum_x \mathbf{1}\{y(x) = 1, \hat{y}(x) = 0\} w(x)}{\sum_x \mathbf{1}\{y(x) = 1\} w(x)}, \quad \text{R-FP} = \frac{\sum_x \mathbf{1}\{y(x) = 0, \hat{y}(x) = 1\} w(x)}{\sum_x \mathbf{1}\{y(x) = 0\} w(x)},$$

82 and the Safety Impact Score

$$\text{SIS}_\lambda = \lambda \text{R-FN} + (1 - \lambda) \text{R-FP}, \quad \lambda \in [0, 1].$$

83 SIS is bounded in $[0, 1]$, interpretable across cases, and exposes FN versus FP asymmetry through λ .
 84 This differs from proximity-weighted Dice, which remains an overlap score and does not separate
 85 FN from FP or normalize by total risk McCullum et al. (2024); Ni et al. (2025).

86 **Hazard-aware tail emphasis.** Let $E_{\text{FN}} = \{x : y = 1, \hat{y} = 0\}$, $E_{\text{FP}} = \{x : y = 0, \hat{y} = 1\}$,
 87 and $h(x) = w(x)$. For $\alpha \in (0, 1]$, let $q_{1-\alpha}$ be the $(1 - \alpha)$ quantile of $h(x)$ over a set E and define
 88 $\text{CVaR}_\alpha(h \mid E) = \frac{1}{\alpha|E|} \sum_{x \in E} h(x) \mathbf{1}\{h(x) \geq q_{1-\alpha}\}$. The Safety Tail Risk is

$$\text{STAR}_{\lambda, \alpha} = \lambda \text{CVaR}_\alpha(h \mid E_{\text{FN}}) + (1 - \lambda) \text{CVaR}_\alpha(h \mid E_{\text{FP}}).$$

89 STAR targets rare but catastrophic mistakes near protected anatomy, something boundary percentiles
 90 and region heuristics do not directly capture Cheng et al. (2021); Huttenlocher et al. (1993).

91 **Theory and invariants.** Because SIS and STAR integrate over a bounded hazard measure in
 92 millimeters, they are physically meaningful, invariant to voxel anisotropy when distances are in
 93 millimeters, and monotone when errors are moved closer to hazard. With a uniform hazard $w(x) \equiv 1$,
 94 SIS reduces to a standard weighted error rate, which shows that our construction strictly generalizes
 95 common accuracy measures.

96 **Reproducible evidence.** We provide a matched-Dice stress protocol and open-dataset scripts that
 97 move equal numbers of boundary voxels toward or away from hazard. Across three MSD tasks,
 98 SIS and especially STAR separate risky from neutral variants while Dice remains unchanged and
 99 hazard-weighted Dice moves little. This complements outcome-focused frameworks by offering a
 100 fast, anatomy-aware safety readout that can be run alongside classical metrics Poel et al. (2025).

101 1.5 Relationship to prior art

102 Our use of a distance-derived hazard field connects to safety margins in surgical robotics Dillon
 103 et al. (2016); Siebold et al. (2017). Our risk normalization and FN versus FP trade-off address
 104 long-standing complaints about average scores and symmetric penalties in clinical validation Maier-
 105 Hein et al. (2024); Poel et al. (2021). Our tail operator plays a role similar to worst-case planners
 106 in autonomous systems but applied to perception evaluation rather than control Antonante et al.
 107 (2023). Finally, proximity-weighted Dice becomes a special case of our framework under a binary
 108 ring-shaped hazard and symmetric weighting, which clarifies how our metrics generalize and strictly
 109 extend that line of work McCullum et al. (2024); Ni et al. (2025).

110 **Summary.** Prior work has weighted overlap by importance, emphasized boundary tails, and
 111 evaluated clinical impact, but there is still no bounded, hazard-aware metric that (i) computes
 112 risk-normalized FN and FP shares with explicit asymmetry and (ii) aggregates the worst errors by
 113 *hazard* rather than by distance or region size. SIS and STAR address these gaps with a continuous
 114 anatomy-derived hazard map, clear interpretation, simple properties, and reproducible demonstrations
 115 on public datasets.

116 2 Methodology

117 2.1 Problem setting and notation

118 Let $\Omega \subset \mathbb{Z}^3$ be a voxel grid with physical spacing (s_x, s_y, s_z) in millimeters and voxel volume
 119 $v = s_x s_y s_z$. Let $y : \Omega \rightarrow \{0, 1\}$ be the ground truth mask of a target structure and $\hat{y} : \Omega \rightarrow \{0, 1\}$ a
 120 prediction. Let $\{C_i\}_{i=1}^M$ denote critical structures to be protected, each provided as a binary mask.
 121 Errors are not equal: a false negative near C_i is typically more consequential than one far away. Our
 122 goal is to evaluate \hat{y} with respect to y through a spatially varying hazard field derived from clearance
 123 to $\{C_i\}$.

124 2.2 Anatomy-derived hazard field

125 For each critical structure C_i , we compute the Euclidean distance transform in millimeters, $d_{C_i}(x)$,
 126 using an anisotropy-aware transform with sampling (s_x, s_y, s_z) . We map distance to a unit hazard
 127 via a monotone decay $r : \mathbb{R}_{\geq 0} \rightarrow [0, 1]$:

$$\text{Polynomial clip: } r(d) = \max\{0, 1 - (d/m)^p\}, \quad m > 0, p \geq 1, \quad (1)$$

$$\text{Exponential: } r(d) = \exp(-d/\tau), \quad \tau > 0. \quad (2)$$

128 We combine multiple structures with either a conservative maximum or an additive accumulation
 129 clipped to one:

$$w(x) = \begin{cases} \max_i w_i r(d_{C_i}(x)) & \text{max aggregation,} \\ \min\{1, \sum_i w_i r(d_{C_i}(x))\} & \text{sum aggregation,} \end{cases} \quad (3)$$

130 where $w_i > 0$ are optional per-structure importances. The hazard map $w : \Omega \rightarrow [0, 1]$ reflects
 131 clearance to protected anatomy and serves as a spatial weight for evaluation.

132 2.3 Risk-normalized error shares and Safety Impact Score

133 Define risk-weighted false negative and false positive shares by normalizing the hazard mass of each
 134 error set by the total available hazard mass in the corresponding region:

$$\text{R-FN} = \frac{\sum_{x \in \Omega} \mathbf{1}\{y(x) = 1, \hat{y}(x) = 0\} w(x)}{\sum_{x \in \Omega} \mathbf{1}\{y(x) = 1\} w(x) + \varepsilon}, \quad (4)$$

$$\text{R-FP} = \frac{\sum_{x \in \Omega} \mathbf{1}\{y(x) = 0, \hat{y}(x) = 1\} w(x)}{\sum_{x \in \Omega} \mathbf{1}\{y(x) = 0\} w(x) + \varepsilon}, \quad (5)$$

135 with a small ε for numerical stability. We mix them with a tunable asymmetry $\lambda \in [0, 1]$:

$$\text{SIS}_\lambda = \lambda \text{R-FN} + (1 - \lambda) \text{R-FP}. \quad (6)$$

136 $\text{SIS}_\lambda \in [0, 1]$ is interpretable as the fraction of hazard mass that is misclassified, with explicit control
 137 of the false negative versus false positive trade off.

138 2.4 Hazard-aware tail emphasis and Safety Tail Risk

139 Let $E_{\text{FN}} = \{x \in \Omega : y = 1, \hat{y} = 0\}$ and $E_{\text{FP}} = \{x \in \Omega : y = 0, \hat{y} = 1\}$. Consider the per-error
 140 hazard values $h(x) = w(x) \in [0, 1]$. For $\alpha \in (0, 1]$, define the conditional value at risk over a set E :

$$\text{CVaR}_\alpha(h \mid E) = \frac{1}{\alpha |E|} \sum_{x \in E} h(x) \mathbf{1}\{h(x) \geq q_{1-\alpha}(E)\}, \quad (7)$$

141 where $q_{1-\alpha}(E)$ is the $(1 - \alpha)$ quantile of $\{h(x) : x \in E\}$. The Safety Tail Risk aggregates the worst
 142 fraction of error hazards for both error types:

$$\text{STAR}_{\lambda, \alpha} = \lambda \text{CVaR}_\alpha(h \mid E_{\text{FN}}) + (1 - \lambda) \text{CVaR}_\alpha(h \mid E_{\text{FP}}). \quad (8)$$

143 $\text{STAR}_{\lambda, \alpha} \in [0, 1]$ focuses on rare but catastrophic mistakes near protected anatomy.

144 2.5 Properties

145 **Boundedness.** Since $h \in [0, 1]$ and normalizers are strictly positive, $\text{R-FN}, \text{R-FP} \in [0, 1]$ and
 146 $\text{SIS}_\lambda \in [0, 1]$. Similarly, $\text{STAR}_{\lambda, \alpha} \in [0, 1]$.

147 **Monotonicity.** Moving any error voxel to a location with larger hazard $w(x)$ weakly increases
 148 SIS_λ and $\text{STAR}_{\lambda, \alpha}$, with denominators fixed.

149 **Reduction.** With $w \equiv 1$, SIS_λ reduces to a standard weighted error rate that generalizes voxel
 150 accuracy. STAR reduces to a tail average over uniform hazard.

151 **Physical invariance.** When distances are computed in millimeters, results are invariant to voxel
 152 anisotropy.

2.6 Baselines and comparators

We report standard and safety-aware baselines:

- Dice coefficient and HD95 as classical overlap and tail-boundary baselines.
- Hazard-weighted Dice (wDice): define

$$\text{wDice} = \frac{2 \sum_x w(x) y(x) \hat{y}(x)}{\sum_x w(x) y(x) + \sum_x w(x) \hat{y}(x) + \varepsilon},$$

which reweights overlap by w for comparison with overlap-style metrics.

- Proximity-weighted Dice variants in radiotherapy when applicable, treated as external baselines.

2.7 Matched-Dice stress protocol

To demonstrate added value beyond overlap, we design a model-free stress test that holds Dice approximately constant while relocating the same number of boundary voxels toward versus away from hazard.

1. Given a ground truth mask y and hazard map w , compute the inner boundary voxels B_{in} and outer boundary voxels B_{out} using a one-voxel erosion or morphological gradient.
2. Rank $B_{\text{in}} \cap \{y = 1\}$ and $B_{\text{out}} \cap \{y = 0\}$ by $w(x)$.
3. For a chosen k , construct two predictions:
 - risk-seeking: flip k highest-hazard inner boundary voxels to false negatives and k highest-hazard outer boundary voxels to false positives,
 - risk-neutral: flip k lowest-hazard inner boundary voxels to false negatives and k lowest-hazard outer boundary voxels to false positives.
4. Report Dice, HD95, wDice, SIS $_{\lambda}$, and STAR $_{\lambda, \alpha}$ for both variants.

This protocol forces classical scores to move minimally while safety-aware metrics diverge when errors hug the hazard.

2.8 Datasets and label mappings

We use three public MSD tasks for generalization:

- **Task08 Hepatic Vessel** target is tumor (label 2) and hazard is vessels (label 1).
- **Task03 Liver** target is tumor (label 2) and hazard is liver (label 1).
- **Task07 Pancreas** target is tumor (label 2) and hazard is pancreas (label 1).

For each case, we read spacing from NIFTI headers to compute millimeter distances. We construct w with a polynomial kernel by default with $m = 10$ mm and $p = 2$, and also report an exponential kernel sensitivity with $\tau = 8$ mm.

2.9 Implementation details

Distances and morphology. We compute $d_{C_i}(x)$ via an anisotropy-aware Euclidean distance transform with sampling (s_x, s_y, s_z) . Boundary sets use a 3D ball structuring element of radius 1 voxel. We ensure the same number of flipped voxels for risky and neutral variants at each case.

Numerics. All sums are carried out in 32-bit float. We use $\varepsilon = 10^{-8}$ in denominators. Quantiles for CVaR are computed with linear interpolation. We verify invariance to regridding by resampling select cases to isotropic spacing and confirming metrics change within tolerance.

Hyperparameters. Unless stated otherwise, we set $\lambda = 0.7$ to overweight false negatives and $\alpha = 0.05$ to summarize the top 5 percent most hazardous mistakes. We ablate $m \in \{5, 10, 15, 20\}$ mm, $p \in \{1, 2, 4\}$, $\tau \in \{6, 8, 10\}$ mm, and $\alpha \in \{0.01, 0.05, 0.1\}$.

Dataset	n	mean Δ SIS \pm CI	p (SIS)	mean Δ STAR \pm CI	p (STAR)	mean Δ wDice
Liver	5	0.008742 ± 0.013258	0.03125	0.368505 ± 0.129822	0.03125	-0.087578
Pancreas	5	0.012759 ± 0.009152	0.03125	0.310152 ± 0.010858	0.03125	-0.117252
H. Vessel	5	0.223996 ± 0.181814	0.03125	0.614249 ± 0.307434	0.03125	-0.421403
Combined	15	0.0818 ± 0.0771	3.052×10^{-5}	0.4310 ± 0.1240	3.242×10^{-4}	–

Table 1: Effect sizes and significance for the matched–Dice stress test (polynomial kernel). CI denotes 95 percent confidence interval on the mean. p values from paired Wilcoxon signed–rank tests with alternative risky $>$ neutral. wDice reported as mean Δ across cases.

193 **Multiple hazards.** For multiple OARs we use max aggregation by default as a conservative choice
194 and report sum aggregation as a sensitivity analysis.

195 2.10 Statistical analysis

196 We perform paired tests between risky and neutral variants within each dataset. For each metric
197 we compute per-case deltas and report mean with 95 percent confidence intervals. We use a paired
198 Wilcoxon signed–rank test with one-sided alternative that risky exceeds neutral for SIS and STAR. We
199 aggregate across datasets and also report per-dataset results. For model-based settings, we compare
200 methods across seeds and report bootstrap confidence intervals where appropriate.

201 3 Results

202 3.1 Experimental setup recap

203 We evaluate the proposed metrics using a matched–Dice stress protocol that holds overlap approxi-
204 mately constant while relocating equal numbers of boundary voxels toward versus away from hazard.
205 We run on three Medical Segmentation Decathlon tasks with five cases each: Hepatic Vessel (target
206 tumor, hazard vessels), Liver (target tumor, hazard liver), and Pancreas (target tumor, hazard pancreas).
207 Unless stated otherwise we use a polynomial hazard kernel with margin $m=10$ mm and exponent
208 $p=2$, tail fraction $\alpha=0.05$, and asymmetry $\lambda=0.7$.

209 3.2 Primary finding: safety metrics separate risky from neutral while Dice does not

210 Across all datasets the tail metric Δ STAR \equiv STAR_{risky}–STAR_{neutral} is strongly positive and
211 statistically significant, indicating that equal–count errors placed nearer to hazard are penalized far
212 more than those pushed away. The risk–normalized Δ SIS is also positive for all datasets, with the
213 largest separation on Hepatic Vessel. By design Δ Dice is zero and provides no safety discrimination,
214 so we do not plot it.

215 3.3 Overlap baselines move little even when risk changes

216 We compare against a hazard–weighted Dice baseline (wDice), which reweights overlap by the hazard
217 field. Although wDice shows some signal, it is substantially less consistent and smaller in magnitude
218 than STAR.

219 3.4 Quantitative summary

220 Table 1 reports effect sizes and paired tests for the polynomial kernel. Hepatic Vessel shows the
221 largest separation, consistent with the clinical intuition that tumors approaching vessels carry high
222 risk. Combined across datasets, both Δ SIS and Δ STAR are significant.

223 For completeness, Table 2 lists per–dataset means of risky and neutral metrics. STAR shows large
224 absolute differences even when Dice is unchanged.

225 3.5 Sensitivity to hazard kernel choice

226 We repeat the analysis with an exponential kernel ($\tau=8$ mm). Direction and relative magnitudes are
227 preserved, indicating robustness to the kernel parameterization.

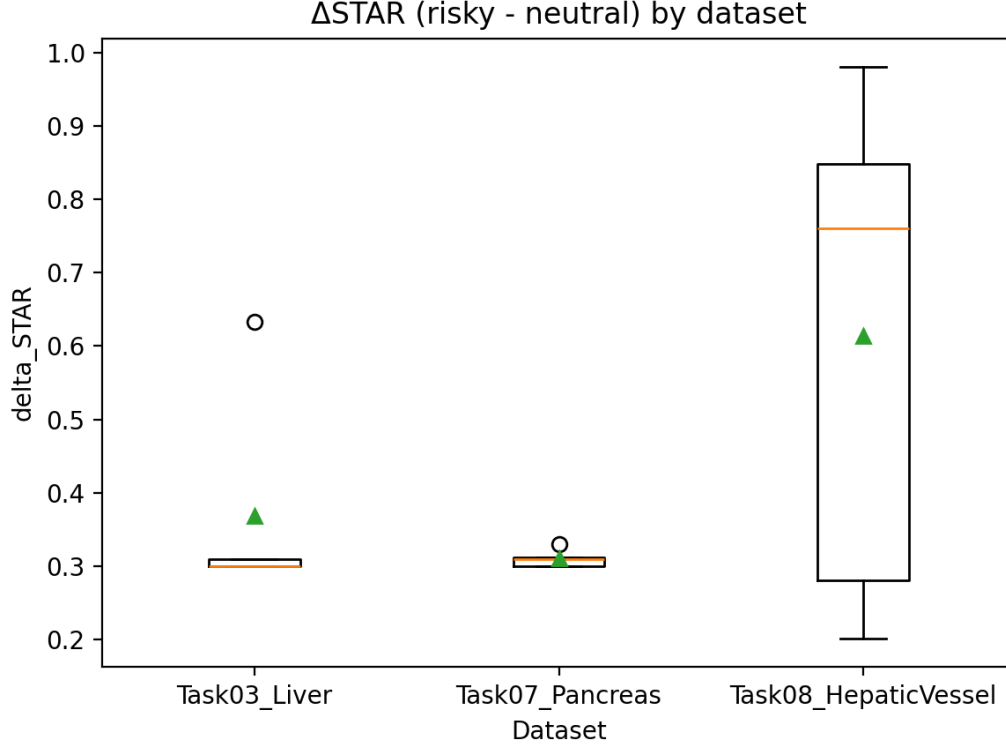


Figure 1: Per-dataset boxplots of ΔSTAR (risky minus neutral). Larger is worse for safety. All three datasets show a consistent positive gap, with the largest effect on Hepatic Vessel where vessels are an explicit organ at risk.

Dataset	Dice _r	Dice _n	SIS _r	SIS _n	STAR _r	STAR _n
Task03 Liver	0.6884	0.6884	0.2293	0.2205	0.9965	0.6280
Task07 Pancreas	0.7744	0.7744	0.1970	0.1843	0.9953	0.6852
Task08 HepaticVessel	0.6556	0.6556	0.4205	0.1965	0.8443	0.2300

Table 2: Per-dataset means for risky (r) and neutral (n) variants under the polynomial kernel. Dice is unchanged by construction while SIS and STAR separate.

3.6 Qualitative exemplars

To illustrate what the hazard field encodes, Figure 4 shows the hazard maps for Hepatic Vessel and Pancreas. These fields convert millimeter clearance to a spatial risk prior that aligns with clinical intuition and explains the observed metric behavior. We omit Task03 Liver exemplars since they are uninformative in this setup.

3.7 Observed edge cases and interpretation

On several Hepatic Vessel cases the neutral variant attains STAR near zero under the clipped polynomial kernel because neutral flips are pushed beyond the 10 mm margin where hazard weights clip to zero. This behavior is expected with a hard margin. The exponential kernel removes clipping and yields nonzero neutral STAR while preserving direction and magnitude, as in Table 3.

Summary SIS and STAR consistently penalize error relocations toward hazard across three datasets while Dice remains flat and wDice moves little. The effect is largest when the hazard corresponds to true organs at risk, and it persists under alternate hazard kernels. These results support the thesis that risk-normalized and tail-focused evaluation captures safety-relevant differences that overlap metrics miss.

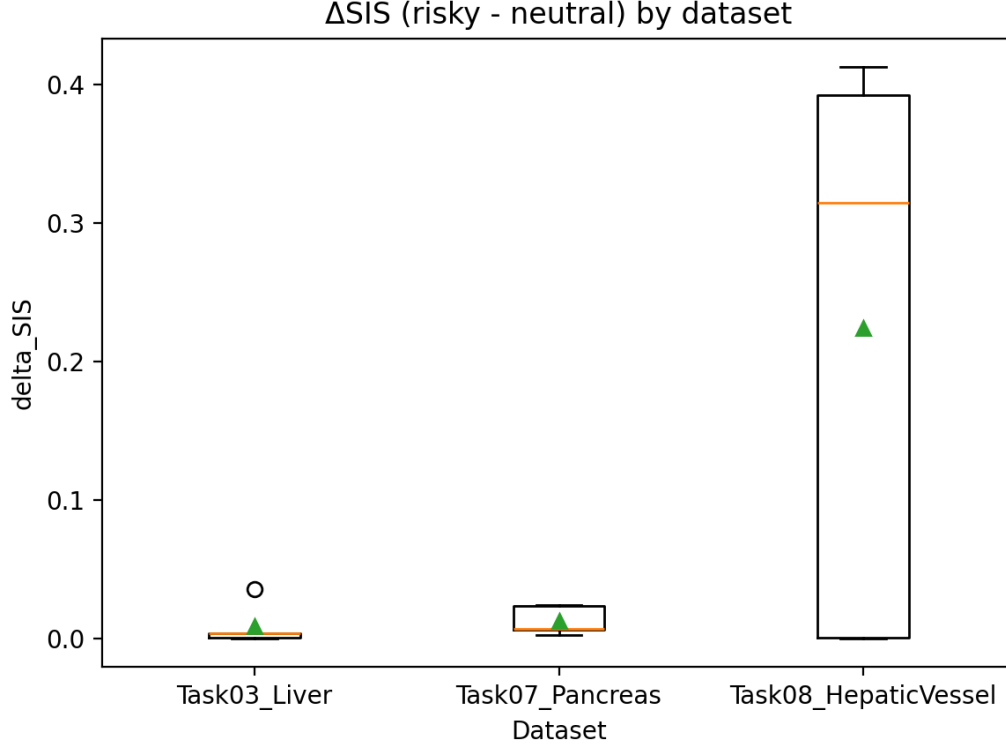


Figure 2: Per-dataset boxplots of Δ SIS (risky minus neutral). SIS separates risky from neutral on all datasets, again with the largest gap on Hepatic Vessel.

Dataset	mean Δ SIS (exp)	mean Δ STAR (exp)
Task03 Liver	0.017233	0.379515
Task07 Pancreas	0.016671	0.341277
Task08 HepaticVessel	0.092108	0.472684

Table 3: Kernel sensitivity with exponential hazard. Effects remain positive and of similar order.

References

- Pasquale Antonante, Sushant Veer, Karen Leung, Xinshuo Weng, Luca Carlone, and Marco Pavone. Task-aware risk estimation of perception failures for autonomous vehicles. In *Robotics: Science and Systems (RSS)*, 2023. URL <https://www.roboticsproceedings.org/rss19/p100.pdf>.
- Chih-Hong Cheng, Alois Knoll, and Hsuan-Cheng Liao. Safety metrics for semantic segmentation in autonomous driving. *arXiv preprint arXiv:2105.10142*, 2021. URL <https://arxiv.org/abs/2105.10142>.
- Neal P. Dillon, Michael A. Siebold, et al. Increasing safety of a robotic system for inner ear surgery using probabilistic error modeling near vital anatomy. In *Proc. SPIE Medical Imaging*, 2016. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC5708556/>.
- Da He, Jayaram K. Udupa, Yubing Tong, and Drew A. Torigian. Mendability index: A new metric for estimating the effort required for manually editing auto-segmentations of objects of interest. In *Proc. SPIE Medical Imaging*, volume 12469, pp. 1246905, 2023. doi: 10.1117/12.2654421. URL <https://pubmed.ncbi.nlm.nih.gov/37256076/>.
- Da He et al. Predicting the effort required to manually mend auto-segmentations in radiotherapy. *medRxiv*, 2024. URL <https://www.medrxiv.org/content/10.1101/2024.06.12.24308779v1>. Analysis of APL, surface Dice, and Mendability Index vs. human effort.

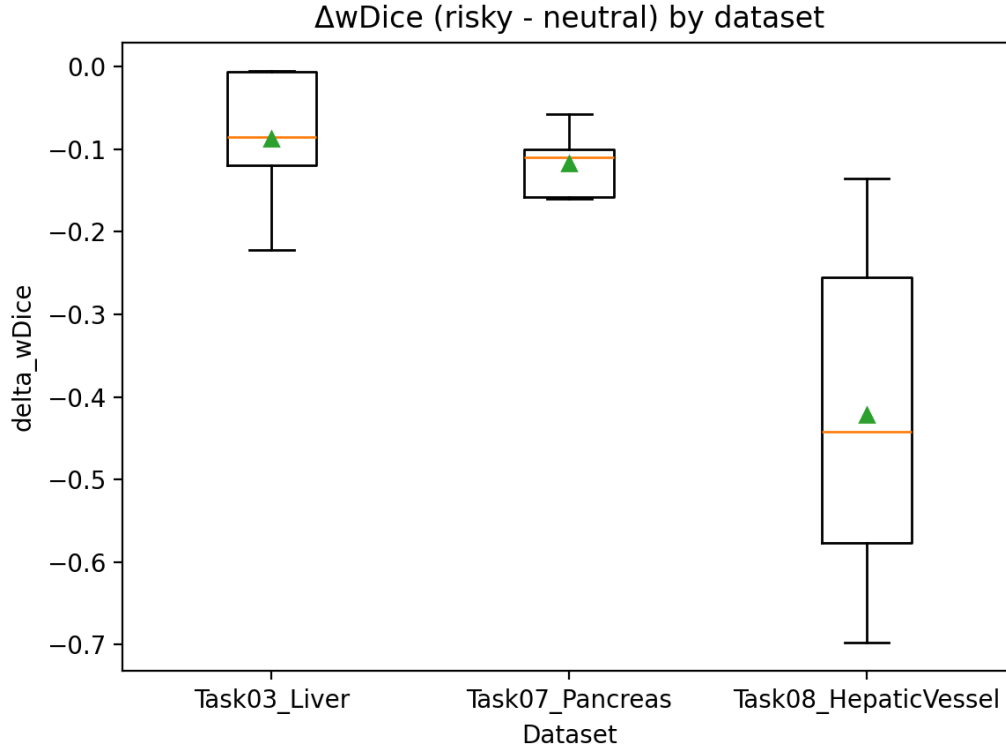


Figure 3: Per-dataset boxplots of $\Delta wDice$ (risky minus neutral). Values are close to zero with high variability compared to $\Delta STAR$. Negative means the risky variant scores worse under wDice.

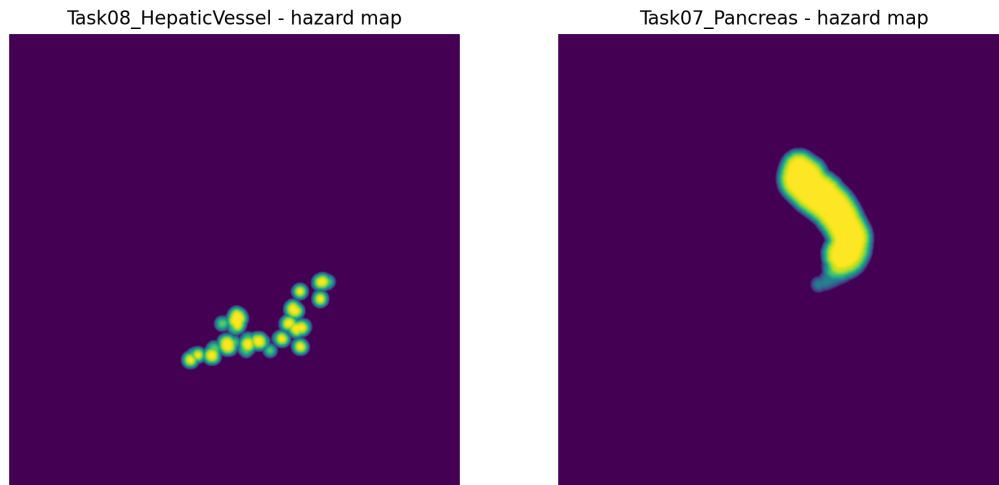


Figure 4: Hazard maps derived from distance to protected anatomy. Left: Hepatic Vessel (vessels as hazard). Right: Pancreas (organ as hazard). Brighter indicates higher hazard near critical structures.

- 260 Daniel P. Huttenlocher, Gregory A. Klanderman, and William J. Rucklidge. Comparing images using
 261 the hausdorff distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(9):850–
 262 863, 1993. doi: 10.1109/34.232073. URL <https://dl.acm.org/doi/10.1109/34.232073>.
- 263 Davood Karimi and Septimiu E. Salcudean. Reducing the hausdorff distance in medical image
 264 segmentation with convolutional neural networks. *arXiv preprint arXiv:1904.10030*, 2019. URL
 265 <https://arxiv.org/abs/1904.10030>.

266 Jun Ma et al. Segment anything in medical images. *Nature Communications*, 15(1160), 2024. URL
267 <https://www.nature.com/articles/s41467-024-44824-z>. Uses NSD alongside DSC per
268 Metrics Reloaded guidance.

269 Lena Maier-Hein, Annika Reinke, Patrick Godau, et al. Metrics reloaded: recommendations for image
270 analysis validation. *Nature Methods*, 21(2):195–212, 2024. doi: 10.1038/s41592-023-02151-z.
271 URL <https://www.nature.com/articles/s41592-023-02151-z>.

272 Lucas McCullum, Kareem A. Wahid, Barbara Marquez, and Clifton D. Fuller. Oar-weighted
273 dice score: A spatially aware, radiosensitivity aware metric for target structure contour quality
274 assessment. *arXiv preprint arXiv:2410.20243*, 2024. URL [https://arxiv.org/abs/2410.](https://arxiv.org/abs/2410.20243)
275 20243.

276 Rui Ni et al. Geometrically focused training and evaluation of organs-at-risk segmentation via
277 deep learning. *Medical Physics*, 2025. URL [https://pmc.ncbi.nlm.nih.gov/articles/](https://pmc.ncbi.nlm.nih.gov/articles/PMC12257911/)
278 PMC12257911/. Proposes weighted Dice (wDSC) using distance to high-risk CTV; reports stronger
279 correlation with dosimetry.

280 Rianne Poel et al. The predictive value of segmentation metrics on dosimetry in organs at risk of
281 the brain. *Medical Image Analysis*, 73:102161, 2021. doi: 10.1016/j.media.2021.102161. URL
282 <https://pubmed.ncbi.nlm.nih.gov/34293536/>.

283 Rianne Poel et al. A comprehensive multifaceted technical evaluation framework for auto-
284 segmentation in radiotherapy. *Communications Medicine*, 2025. URL [https://www.nature.](https://www.nature.com/articles/s43856-025-01048-6)
285 com/articles/s43856-025-01048-6. Evaluation across geometry, efficiency, treatment qual-
286 ity, and expert acceptability.

287 Annika Reinke et al. Understanding metric-related pitfalls in image analysis validation. *arXiv*
288 *preprint arXiv:2302.01790*, 2023. URL <https://arxiv.org/abs/2302.01790>.

289 Michael A. Siebold, Loris Fichera, Neal P. Dillon, and J. Michael Fitzpatrick. Safety margins in
290 robotic bone milling: from registration uncertainty to statistically safe surgeries. *The International*
291 *Journal of Medical Robotics and Computer Assisted Surgery*, 13(4):e1773, 2017. doi: 10.1002/rcs.
292 1773. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/rcs.1773>.

293 Various. Groundbreaking insights into the implementation of metrics in medical image segmentation.
294 *arXiv preprint arXiv:2410.02630*, 2024. URL <https://arxiv.org/abs/2410.02630>. Contains
295 discussion of Normalized Surface Dice.

296 Y. Zhang et al. Comprehensive clinical usability-oriented contour quality assessment in radiotherapy.
297 *Physica Medica (or similar clinical engineering journal)*, 2024. URL [https://pmc.ncbi.](https://pmc.ncbi.nlm.nih.gov/articles/PMC11711007/)
298 nlm.nih.gov/articles/PMC11711007/. Discusses correlation of APL and Surface Dice with
299 contour efficiency.