

Controlling False Positives in Image Segmentation via Conformal Prediction

Luca Mossina Corentin Friedrich

IRT Saint Exupéry, Toulouse, France

Abstract

Reliable semantic segmentation is essential for clinical decision making, yet deep models rarely provide explicit statistical guarantees on their errors. We introduce a simple post-hoc framework that constructs confidence masks with distribution-free, image-level control of false-positive predictions. Given any pretrained segmentation model, we define a nested family of shrunk masks obtained either by increasing the score threshold or by applying morphological erosion. A labeled calibration set is used to select a single shrink parameter via conformal prediction, ensuring that, for new images that are exchangeable with the calibration data, the proportion of false positives retained in the confidence mask stays below a user-specified tolerance with high probability. The method is model-agnostic, requires no retraining, and provides finite-sample guarantees regardless of the underlying predictor. Experiments on a polyp-segmentation benchmark demonstrate target-level empirical validity. Our framework enables practical, risk-aware segmentation in settings where over-segmentation can have clinical consequences. Code at <https://github.com/deel-ai-papers/conseco>.

1 Introduction

Reliable segmentation is a prerequisite for clinical use of deep-learning models, where false positives may trigger unnecessary interventions. Existing uncertainty scores and calibration methods provide useful heuristics, but they do not offer finite-sample guarantees on the errors of the produced masks. Given any pretrained segmentation model, our post-hoc method builds inner masks $I_\lambda(X)$ by progressively shrinking the predicted mask \hat{Y} using a single control parameter λ , through either sigmoid score thresholding or morphological erosion. We calibrate the shrinkage level on a small held-out labeled set.

Our procedure, based on inductive (or “split”) Conformal Prediction (CP) [21, 15], guarantees that at a user-chosen confidence level $1 - \alpha$, the inner mask contains at most a user-specified fraction τ of false-positive pixels. The validity bound is asserted at the *image level*; we refer to the inner mask $I_\lambda(X)$ as “confidence mask” and the remainder of the prediction is flagged as uncertain, producing the uncertainty region $U_\lambda(X) = \hat{Y} \setminus I_\lambda(X)$.

Contributions. (i) A black-box, distribution-free method that returns inner prediction sets I_λ with guaranteed control of accepted false positives at the image level. (ii) Two concrete, implementation-ready formulations: score thresholding and morphological erosion. (iii) A calibration protocol that exposes two operational choices to the user, confidence $1 - \alpha$

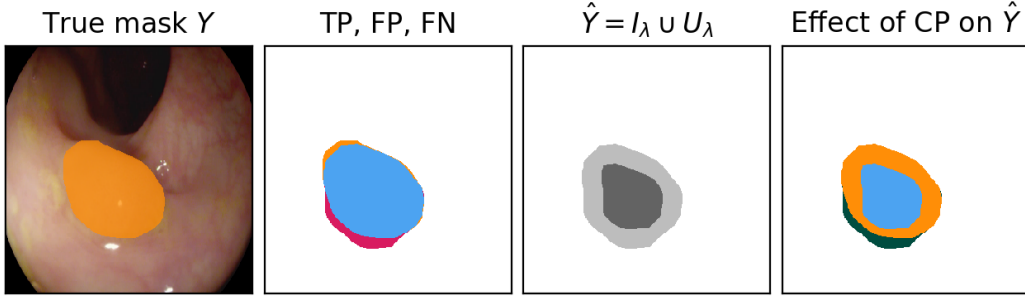


Figure 1: Example with erosion inner mask $I_\lambda^\varepsilon(X)$ at $\tau = 0.01$ and $1 - \alpha = 0.9$. From the left: (i) Ground-truth mask Y overlayed on input image X ; (ii) true positives & false positives in \hat{Y} , and Y pixels missed (false negatives); (iii) Inner “confidence” mask $I_\lambda^\varepsilon(X)$ (dark grey) and uncertainty “rejection” mask $U_\lambda(X)$ (light grey); (iv) \hat{Y} is shrunk to $I_\lambda^\varepsilon(X)$. $U_\lambda(X)$ rejects most FPs (■) but also some TPs, i.e. g.t. pixels (■) well-segmented in \hat{Y} . **Colors.** ■: true mask Y ; ■: false positives (FP); ■: true positives (TP); ■: rejected FPs.

and tolerated accepted false positives τ , with no retraining. (iv) Evidence on a biomedical benchmark that the empirical image-level validity of the *accepted false-positive* proportion (AFP) matches the target confidence.

2 Related Work

Split CP [15] constructs distribution-free prediction sets with finite-sample guarantees of containing the true target at a user-specified confidence level $1 - \alpha$. Conformal Risk Control extends CP to monotone losses, providing guarantees on the expected risk [2]. In multilabel prediction, inner and outer sets are constructed to enclose the true label set [6].

For CP in semantic segmentation, [8] have used inner and outer prediction masks targeting coverage of the ground-truth mask. Other conformal approaches reduce false negatives by lowering score thresholds [2, 13, 5] or by morphological dilation of the predicted mask [14]. CP is complementary to the broader literature on uncertainty quantification: methods such as MC-Dropout [11], deep ensembles [12], or failure prediction [7] provide useful uncertainty maps but do not yield distribution-free guarantees at deployment, which can be achieved with CP.

Our positioning. We address binary medical segmentation and control the proportion of accepted false-positive pixels within the predicted region using nested prediction-shrinkers and conformal calibration. This is complementary to prior false-negative control work and different in scope from multilabel FP-limited set prediction [10]. We instantiate the nested sets through standard anti-extensive morphological operators [17] and sigmoid-score thresholding.

3 Methods

We aim to control statistically the number of false-positive pixels accepted in predicted masks. To do so, we define a quantity $\mathcal{F}_\lambda(X, Y)$ (Eq. 3) that is compatible with the requirements of CP and hence admits rigorous statistical guarantees, notably via the inner sets proposed in Eq. 1 and Eq. 2. We refer to \mathcal{F}_λ as the *accepted false-positive proportion (AFP)*: it measures the fraction of false positives that remain in the accepted region relative to the original predicted area $|\hat{Y}|$.

Let X be an image over a grid $\Omega \subseteq \mathbb{Z}^2$ of $n_H \times n_W$ pixels, and let $Y \subseteq \Omega$ and $\hat{Y} \subseteq \Omega$ denote the ground-truth and predicted segmentation masks obtained with a segmentation model, respectively.

Defining inner prediction sets. We construct a nested family of *inner prediction sets* $\{I_\lambda(X)\}_{\lambda \in \Lambda}$ such that: (i) $I_\lambda(X) \subseteq \hat{Y}$, (ii) there exists λ_0 such that $I_{\lambda_0}(X) = \hat{Y}$, and (iii) for any $\lambda_1 \leq \lambda_2$, $I_{\lambda_1}(X) \supseteq I_{\lambda_2}(X)$. We interpret the inner prediction sets as **confidence masks** at a chosen confidence level, defining subregions of the prediction that are “accepted” according to the conformal procedure.

We propose two simple ways to shrink a predicted mask so that fewer false-positive pixels are accepted as confident. We restrict our exposition to two inner set models that are applicable *a posteriori* to most segmentation models, although any nested family of sets (see above) can be used.

First is a set that applies a **threshold on sigmoid** scores $\hat{\sigma}(X)_{ij}$, the output of a binary segmentation model:

$$I_\lambda^\sigma(X) := \{\text{pixels } (i, j) \text{ s.t. } \hat{\sigma}(X)_{ij} \geq \lambda\}, \quad (1)$$

with $\lambda \in [0.5, 1]$, where $\lambda_0 = 0.5$ is the threshold commonly used in segmentation.

Second, we build a set that works as a dual to the morphological dilation [17] used in [14]: we apply **morphological erosion** $\varepsilon_B(\cdot)$ to the mask \hat{Y} as

$$I_\lambda^\varepsilon(X) := \underbrace{(\varepsilon_B \circ \varepsilon_B \circ \dots \circ \varepsilon_B)}_{\lambda \text{ iterations}}(\hat{Y}) = \varepsilon_B^\lambda(\hat{Y}). \quad (2)$$

We fix a structuring element, e.g. $B = \begin{bmatrix} \blacksquare & \blacksquare \\ \blacksquare & \blacksquare \end{bmatrix}$ for 4-connectivity, and erode \hat{Y} λ times, $\lambda \in \mathbb{N}$; note that for $\lambda_0 = 0$, $I_{\lambda_0}(X) = \hat{Y}$. This erosion model also applies to black-box predictors, whose internals are not accessible to end users (e.g. third-party vendors, embedded in medical equipment, etc.). This function is suited to segmentation models where the false positives are concentrated at the boundary of the object. As noted in [14], *any* morphology-based inner set is applicable, e.g., combining structuring elements or using a discretized ball whose radius is controlled by λ .

Our inner sets are nested. Our definitions give rise to nested sets, that is, for any $\lambda_1 \leq \lambda_2$, we get $I_{\lambda_1} \supseteq I_{\lambda_2}$. For $I_\lambda^\sigma(X)$, as λ grows, fewer pixels in \hat{Y} have scores above this threshold and fewer pixels are included in I_λ^σ (Eq. 1), hence its size is non-increasing in λ . For I_λ^ε , since morphological erosion is anti-extensive (i.e. contractive), we have $\varepsilon_B^{\lambda_1}(\hat{Y}) \supseteq \varepsilon_B^{\lambda_2}(\hat{Y})$ for $\lambda_1 \leq \lambda_2$.

3.1 Formulation of the False-Positives Control problem

Inner sets $I_\lambda(X)$ aim to ignore false positives in predictions. However, due to noise in predictive models or annotation errors in segmentation datasets, reaching zero false-positive pixels

would require large values of λ . The obtained inner sets would thus be very small or even empty; note that the trivial solution $I_\lambda = \emptyset$ does not contain any FPs and it is always valid. To avoid this trivial solution, we allow a small fraction of FPs, which is controlled by a user-defined tolerance parameter $\tau \in [0, 1]$.

Let $W(X, Y) = \hat{Y} \cap (\Omega \setminus Y)$ denote the set of **false-positive** (FP) pixels in \hat{Y} . This set does not depend on the inner mask I_λ . Since we cannot control the size of $W(X, Y)$, which depends on the fixed predictor, we instead construct inner masks that are the largest subsets of \hat{Y} containing few false positives, according to the chosen definition of inner mask (Eqs. 1–2). We refer to $I_\lambda(X)$ as the *confidence mask* and to its complement $U_\lambda(X) = \hat{Y} \setminus I_\lambda(X)$ as the *uncertain region*.

We want to control the following quantity, representing the **accepted false-positive proportion** (AFP) within the predicted mask \hat{Y} :

$$\mathcal{F}_\lambda(X, Y) = \frac{|I_\lambda(X) \cap W(X, Y)|}{|\hat{Y}|}. \quad (3)$$

Here $|\cdot|$ denotes set cardinality. If $|\hat{Y}| = 0$, we set $I_\lambda = \emptyset$ and $\mathcal{F}_\lambda = 0$. Since I_λ is nested and the denominator $|\hat{Y}|$ is λ -invariant, \mathcal{F}_λ is non-increasing in λ . Dividing by $|I_\lambda|$ would break this monotonicity and is therefore avoided.

Interpretation. \mathcal{F}_λ quantifies the proportion of false positives that remain “accepted” within the confidence mask at a level λ . For $\lambda = \lambda_0$ (no shrinkage), $\mathcal{F}_{\lambda_0} = |W(X, Y)|/|\hat{Y}|$ corresponds to the original false-positive fraction of the prediction. Increasing λ enforces stricter acceptance and can only decrease \mathcal{F}_λ . In practice, \mathcal{F}_λ expresses the fraction of spurious detections that remain unfiltered after applying the confidence threshold; in some cases it is possible to have $\mathcal{F}_\lambda(X, Y) \leq \tau$, in which case no shrinkage would be needed.

3.2 Conformal Prediction

Our method builds on the standard inductive CP framework [15]. We adapt the inner sets from the segmentation approach of [8], and use morphological erosion as the counterpart of the dilation-based outer sets in [14]. Importantly, CP provides *marginal frequentist* guarantees on the *mask-level procedure*, rather than on individual pixels: if the calibration and testing process were repeated many times, the empirical validity condition in Eq. (6) would be satisfied in at least $100(1 - \alpha)\%$ of cases on average. This states that, for exchangeable data, the accepted false-positive proportion (AFP) of the inner mask satisfies $\mathcal{F}_{\hat{\lambda}} \leq \tau$ with probability at least $1 - \alpha$ at the image level. Once the user has fixed a tolerance value $\tau \in [0, 1]$, the nonconformity score $r_\iota = r(X_\iota, Y_\iota)$ for a calibration pair (X_ι, Y_ι) is computed as (details in Sec. 3.3)

$$r_\iota = \inf \{ \lambda : \mathcal{F}_\lambda(X_\iota, Y_\iota) \leq \tau \}, \quad (4)$$

and the conformalizing quantile is the

$$\hat{\lambda} = \lceil (n + 1)(1 - \alpha) \rceil\text{-th smallest value in } (r_\iota)_{\iota=1}^n. \quad (5)$$

Assuming that calibration data and test point $(X_{\text{new}}, Y_{\text{new}})$ form an exchangeable (or i.i.d.) sequence, we obtain the distribution-free, marginal guarantee:

$$\mathbb{P}(\mathcal{F}_{\hat{\lambda}}(X_{\text{new}}, Y_{\text{new}}) \leq \tau) \geq 1 - \alpha. \quad (6)$$

Proof. Because inner masks $I_\lambda(X)$ are nested with respect to λ , \mathcal{F}_λ is non-increasing in λ . Define the binary loss $\ell(X, Y, \lambda) = \mathbb{1}\{\mathcal{F}_\lambda(X, Y) > \tau\}$, which is monotone non-increasing in λ . Applying Conformalized Risk Control [2] with this loss, whose CRC selection rule coincides with the conformal quantile (Sec. 2.3 in [2]), yields Eq. (6).

Interpretation. This image-level guarantee states that, with probability (i.e. confidence level) at least $1 - \alpha$, the region of the inner mask $I_{\hat{\lambda}}(X_{\text{new}})$ will not “accept” more than a fraction τ (e.g., 5%) false positives in the original prediction \hat{Y}_{new} . This can also be seen as a mechanism to produce confidence regions that mitigate the operational hazard of predicting a treatment (i.e. a “positive” pixel) where not needed. Clinically, AFP bounds the fraction of unnecessary positive pixels *within the region we agree to act upon*, while the rejected pixels are marked as uncertain and not acted on.

Scope. Our guarantee depends only on exchangeability of the data and on the nestedness of the inner sets $I_\lambda(X)_{\lambda \in \Lambda}$. It does not depend on the particular way I_λ is produced. Any scoring function or post-hoc shrinker that induces a nested family can be plugged in; this choice affects only *utility* (e.g., average contraction, inner-margin size, ATP/CR), not validity. The convention of a fixed denominator $|\hat{Y}|$ and the handling of empty masks ($\mathcal{F}_\lambda = 0$) likewise do not change the guarantee; they only influence how conservative the resulting inner masks are. This opens the method to empirical exploration of shrinkers and scores from any uncertainty model.

3.3 Conformalization procedure

One can use any fixed segmentation model, either pretrained or trained separately on dedicated data.

1. **Collect calibration data.** Gather a labeled calibration set $\{(X_\iota, Y_\iota)\}_{\iota=1}^n$, distinct from both training and test data. Following standard CP assumptions, these samples are required to be exchangeable (or i.i.d.) with the test cases.
2. **Define the family of inner sets.** Select either: (i) the threshold model I_λ^σ in Eq. (1) with $\Lambda = [0.5, 1]$ or (ii) the morphological model I_λ^ε in Eq. (2) with $\Lambda = \mathbb{N}$ and a structuring element, e.g. $B = \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{bmatrix}$ or $\begin{bmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{bmatrix}$.
3. **Compute calibration scores.** For each calibration pair (X_ι, Y_ι) , compute $r_\iota = r(X_\iota, Y_\iota)$ as in Eq. (4).
 - (i) For I_λ^σ , evaluate \mathcal{F}_λ at breakpoints $\lambda \in \{\hat{\sigma}(X_\iota)_{ij} : (i, j) \in \hat{Y}_\iota\} \cup \{1\}$ and keep the smallest λ s.t. $\mathcal{F}_\lambda \leq \tau$.
 - (ii) For I_λ^ε , iterate ε_B until $\mathcal{F}_\lambda \leq \tau$ or the mask becomes empty, in which case $\mathcal{F}_\lambda = 0$.
4. **Estimate the conformal threshold.** Let $k = \lceil (n+1)(1-\alpha) \rceil$, and set $\hat{\lambda}$ to the k -th smallest value among the calibration scores $\{r_\iota\}_{\iota=1}^n$, as in Eq. (5).
5. **Predict on the test set.** For a new input X_{new} , output the inner mask $I_{\hat{\lambda}}(X_{\text{new}})$.

4 Experiments

Following prior work on distribution-free methods for segmentation [3, 2, 5], we evaluate on the POLYPS dataset collection [18, 4, 19, 20, 16] using pretrained PraNet [9].¹ We compare three procedures: *Baseline* (uncalibrated sigmoid threshold at 0.5, i.e., $I_{\lambda_0} = \hat{Y}$), *Threshold* (CP by increasing the score cutoff), and *Erosion* (CP by morphological erosion). Confidence level is $1 - \alpha = 0.90$ and tolerance $\tau \in \{0.1, 0.01, 0.001\}$. For erosion we use a cross structuring element $B = \begin{bmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{bmatrix}$. Each run randomly permutes the test split and assigns half to calibration and half to evaluation, yielding 250 calibration images. We report mean \pm standard deviation over 10 seeds.

Metrics. We quantify validity and utility with three image-level quantities. (i) *Empirical validity* (EV): fraction of test images whose accepted false-positive proportion does not exceed τ , that targets the nominal frequency $1 - \alpha$,

$$\text{EV} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \mathbb{1}\{\mathcal{F}_{\hat{\lambda}}(X_i, Y_i) \leq \tau\}. \quad (7)$$

(ii) *Contraction ratio* (CR): retained area after shrinkage,

$$\text{CR} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \frac{|I_{\hat{\lambda}}(X_i)|}{|\hat{Y}_i|}. \quad (8)$$

CR = 1 for the Baseline and decreases as masks shrink, measuring the utility loss incurred by enforcing statistical validity.

(iii) *Accepted true-positive fraction* (ATP):

$$\text{ATP} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \frac{|I_{\hat{\lambda}}(X_i) \cap Y_i|}{|\hat{Y}_i|}. \quad (9)$$

For the Baseline, $\text{ATP} = \frac{|TP|}{|\hat{Y}|}$ equals image-level precision, i.e. the fraction of TPs in the prediction mask.

4.1 Results

Quantitative results at $1 - \alpha = 0.9$ are reported in Tab. 1. Conformalized procedures attain empirical validity (EV) close to the nominal target across all τ , whereas the Baseline fails for small τ and only approaches validity in the more permissive setting $\tau = 0.1$, confirming the need for calibration. Utility decreases smoothly as τ tightens. At fixed τ , the Threshold variant retains higher CR and ATP than Erosion, reflecting a more moderate shrinkage path on typical polyp masks, while Erosion removes a larger fraction of peripheral predictions. This is expected, as sigmoid scores are often more informative than the binary mask; when available, they should be evaluated alongside erosion. The validity-utility trade-off behaves monotonically and remains stable across random seeds. Finally, since Conformal Prediction guarantees statistical validity for *any* segmentation model, utility metrics can also serve to compare or select among different predictors, e.g. when provided by a third-party.

¹We use the precomputed predictions as distributed by the authors of [3, 2] for their comprehensive introduction to CP [1] at github.com/aangelopoulos/conformal-prediction,

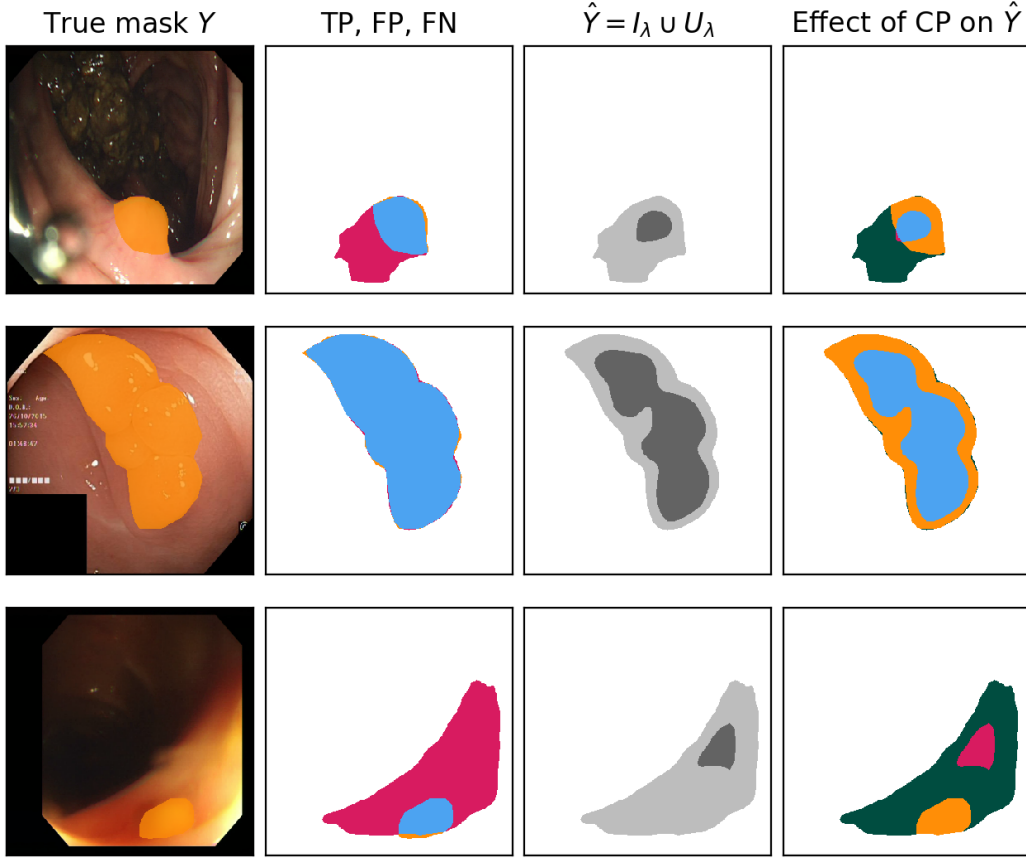


Figure 2: Examples with thresholding inner mask $I_\lambda^\sigma(X)$, at $\tau = 0.01$ and $1 - \alpha = 0.9$. **Top.** Large FP removal with moderate TP loss. **Middle.** When FPs are already negligible, shrinkage removes TPs. **Bottom.** Failure case: residual inner mask is FP-only. **Colors.** ■: true mask Y ; ■: false positives; ■: true positives; ■: rejected false positives.

5 Conclusion and Perspectives

We introduced a post-hoc conformal procedure for binary segmentation that calibrates a single shrink parameter to control, with finite-sample guarantees, the image-level proportion of accepted false positives. The method is model-agnostic, requires no retraining, and operates with either score-thresholding or morphological-erosion shrinkers. Experiments on polyp segmentation demonstrate target-level empirical validity and a smooth validity–utility trade-off governed by τ . Limitations include the mask-level (marginal) nature of the guarantees and potential true-positive loss when predictions are already precise. Future work includes extending the approach to multi-class segmentation and building size- and instance-adaptive inner masks.

	Method	EV	CR (\uparrow)	ATP (\uparrow)
$\tau = 0.1$	Baseline	0.789 ± 0.015	—	0.873 ± 0.014
	Threshold	0.931 ± 0.049	0.701 ± 0.129	0.676 ± 0.111
	Erosion	0.897 ± 0.030	0.578 ± 0.155	0.527 ± 0.132
$\tau = 0.01$	Baseline	0.165 ± 0.013	—	0.873 ± 0.014
	Threshold	0.926 ± 0.037	0.532 ± 0.087	0.525 ± 0.083
	Erosion	0.902 ± 0.034	0.274 ± 0.085	0.252 ± 0.075
$\tau = 0.001$	Baseline	0.005 ± 0.002	—	0.873 ± 0.014
	Threshold	0.914 ± 0.029	0.473 ± 0.067	0.469 ± 0.065
	Erosion	0.911 ± 0.033	0.112 ± 0.052	0.103 ± 0.045

Table 1: Results for confidence level $1 - \alpha = 0.9$.

Acknowledgments

This work was carried out within the DEEL project,² which is part of IRT Saint Exupéry and the ANITI AI cluster. The authors acknowledge the financial support from DEEL’s Industrial and Academic Members and the France 2030 program – Grant agreements n°ANR-10-AIRT-01 and n°ANR-23-IACL-0002.

References

- [1] ANGELOPOULOS, A. N., AND BATES, S. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511* (2021).
- [2] ANGELOPOULOS, A. N., BATES, S., FISCH, A., LEI, L., AND SCHUSTER, T. Conformal risk control. In *The Twelfth International Conference on Learning Representations* (2024).
- [3] BATES, S., ANGELOPOULOS, A., LEI, L., MALIK, J., AND JORDAN, M. Distribution-free, risk-controlling prediction sets. *J. ACM* 68, 6 (9 2021).
- [4] BERNAL, J., SÁNCHEZ, F. J., FERNÁNDEZ-ESPARRACH, G., GIL, D., RODRÍGUEZ, C., AND VILARIÑO, F. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Comp. Med. Imaging Graph.* 43 (2015), 99–111.
- [5] BLOT, V., ANGELOPOULOS, A. N., JORDAN, M., AND BRUNEL, N. J.-B. Automatically adaptive conformal risk control. In *Proceedings of The 28th International Conference on Artificial Intelligence and Statistics* (May 2025), vol. 258 of *PMLR*, pp. 19–27.
- [6] CAUCHOIS, M., GUPTA, S., AND DUCHI, J. C. Knowing what you know: valid and validated confidence sets in multiclass and multilabel prediction. *Journal of Machine Learning Research* 22, 81 (2021), 1–42.

²<https://www.deel.ai>

- [7] CORBIÈRE, C., THOME, N., BAR-HEN, A., CORD, M., AND PÉREZ, P. Addressing failure prediction by learning model confidence. In *Advances in Neural Information Processing Systems* (2019), vol. 32, Curran Associates, Inc.
- [8] DAVENPORT, S. Conformal confidence sets for biomedical image segmentation. *arXiv preprint arXiv:2410.03406* (2024).
- [9] FAN, D.-P., JI, G.-P., ZHOU, T., CHEN, G., FU, H., SHEN, J., AND SHAO, L. Pranet: Parallel reverse attention network for polyp segmentation. In *MICCAI 2020* (2020), Springer, pp. 263–273.
- [10] FISCH, A., SCHUSTER, T., JAAKKOLA, T., AND BARZILAY, D. Conformal prediction sets with limited false positives. In *Proceedings of the 39th International Conference on Machine Learning* (17–23 Jul 2022), vol. 162 of *Proceedings of Machine Learning Research*, PMLR, pp. 6514–6532.
- [11] GAL, Y., AND GHAHRAMANI, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of The 33rd International Conference on Machine Learning* (New York, New York, USA, 20–22 Jun 2016), vol. 48 of *Proceedings of Machine Learning Research*, PMLR, pp. 1050–1059.
- [12] LAKSHMINARAYANAN, B., PRITZEL, A., AND BLUNDELL, C. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems* 30 (2017).
- [13] MOSSINA, L., DALMAU, J., AND ANDÉOL, L. Conformal semantic image segmentation: Post-hoc quantification of predictive uncertainty. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops* (June 2024), pp. 3574–3584.
- [14] MOSSINA, L., AND FRIEDRICH, C. Conformal prediction for image segmentation using morphological prediction sets. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2025* (Cham, 2026), Springer Nature Switzerland, pp. 78–88.
- [15] PAPADOPOULOS, H., PROEDROU, K., VOVK, V., AND GAMMERMAN, A. Inductive confidence machines for regression. In *ECML 2002* (2002).
- [16] POGORELOV, K., RANDEL, K. R., GRIWODZ, C., ESKELAND, S. L., DE LANGE, T., ET AL. Kvasir: A multi-class image dataset for computer aided gastrointestinal disease detection. In *Proceedings of ACM MMSys* (2017), p. 164–169.
- [17] SERRA, J. *Image analysis and mathematical morphology: V.1*. Academic Press, 1984.
- [18] SILVA, J., HISTACE, A., ROMAIN, O., DRAY, X., AND GRANADO, B. Toward embedded detection of polyps in wce images for early diagnosis of colorectal cancer. *Int J CARS* 9 (2014), 283–293.
- [19] TAJBAKHSH, N., GURUDU, S. R., AND LIANG, J. Automated polyp detection in colonoscopy videos using shape and context information. *IEEE T-MI* 35, 2 (2015), 630–644.

- [20] VÁZQUEZ, D., BERNAL, J., SÁNCHEZ, F. J., FERNÁNDEZ-ESPARRACH, G., LÓPEZ, A. M., ROMERO, A., ET AL. A benchmark for endoluminal scene segmentation of colonoscopy images. *J Healthc Eng* 2017, 1 (2017), 4037190.
- [21] VOVK, V., GAMMERMAN, A., AND SHAFER, G. *Algorithmic learning in a random world*, vol. 29. Springer, 2005.