



HAL
open science

Evidential fully convolutional network for semantic segmentation

Zheng Tong, Philippe Xu, Thierry Denoeux

► **To cite this version:**

Zheng Tong, Philippe Xu, Thierry Denoeux. Evidential fully convolutional network for semantic segmentation. *Applied Intelligence*, 2021, 51 (9), pp.6376-6399. 10.1007/s10489-021-02327-0. hal-03511108

HAL Id: hal-03511108

<https://hal.science/hal-03511108>

Submitted on 4 Jan 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Evidential fully convolutional network for semantic segmentation

Zheng Tong · Philippe Xu · Thierry Dencœux

Received: date / Accepted: date

Abstract We propose a hybrid architecture composed of a fully convolutional network (FCN) and a Dempster-Shafer layer for image semantic segmentation. In the so-called evidential FCN (E-FCN), an encoder-decoder architecture first extracts pixel-wise feature maps from an input image. A Dempster-Shafer layer then computes mass functions at each pixel location based on distances to prototypes. Finally, a utility layer performs semantic segmentation from mass functions and allows for imprecise classification of ambiguous pixels and outliers. We propose an end-to-end learning strategy for jointly updating the network parameters, which can make use of soft (imprecise) labels. Experiments using three databases (Pascal VOC 2011, MIT-scene Parsing and SIFT Flow) show that the proposed combination improves the accuracy and calibration of semantic segmentation by assigning confusing pixels to multi-class sets.

Keywords Evidence theory · belief function · fully convolutional network · decision analysis · semantic segmentation

1 Introduction

In the past few decades, one of the most difficult problems in computer vision has been image semantic segmentation, defined as the process of partitioning a digital image into multiple sets of pixels. The result of image segmentation is a set of segments that collectively cover the entire image, called the *segmentation mask*. This mask constitutes a simplified representation, more meaningful and easier to analyze than the original image. Semantic segmentation has been widely applied

Zheng Tong, Philippe Xu
Université de technologie de Compiègne, CNRS, Heudiasyc, Compiègne, France
E-mail: zheng.tong@hds.utc.fr, philippe.xu@hds.utc.fr

Thierry Dencœux
Université de technologie de Compiègne, CNRS, Heudiasyc, Compiègne, France, and
Institut universitaire de France, Paris, France
E-mail: thierry.denoeux@utc.fr

to advanced driver assistance systems [6, 16, 36], human-machine interaction [39], medical imaging [18], and so on.

In the last decade, deep learning-based models, especially fully convolutional networks (FCNs) [25] and variants [22, 28], have been developed for semantic segmentation and have achieved remarkable success. FCNs take advantage of existing deep neural networks, which have the capacity to learn reliable and robust features. An FCN transforms existing and well-known classification models such as VGG (16-layer net) [31] or ResNet [21] into fully convolutional ones by replacing the fully connected layers with convolutional ones to output spatial maps instead of classification scores. Those maps are upsampled using fractionally-strided convolutions (called *deconvolutions* [41, 42]) to produce dense per-pixel labeled outputs. This approach has allowed for significant improvements in segmentation accuracy over traditional methods on benchmark databases like Pascal VOC 2011 [17]. However, despite the power and flexibility of the FCN-based models, they still face the following three problems:

1. *How to perform novelty detection?* In many learning sets, not all classes are labeled, especially for some objects in the background. An ideal image segmentation algorithm should detect “unknown” objects belonging to classes that are not represented in the learning set. This capacity is called *novelty detection* [8]. FCN-based models generally randomly assign unknown objects to one of the known classes, though some models tend to assign unknown objects to the background class.
2. *How to process pixels with confusing information?* In image-segmentation training sets, all pixels are precisely labeled, even if the true label is actually uncertain. This is the case, for example, for the pixels at object borders. Pixels with precise but incorrect labels may have negative effects on the performance of learning systems [2, 27].
3. *When will the FCN-based methods fail?* In decision-making systems, a neural network should not only be as accurate as possible, but it should also have the ability to signal when it is likely to be incorrect. Neural networks developed nowadays tend not to be well calibrated [20], though they are more accurate than they were a decade ago. In other words, the accuracy of modern neural networks, including FCN-based models, does not match their confidence.

In this paper, the above problems are tackled using the Dempster-Shafer (DS) theory of belief functions [7, 30]. This mathematical framework, also referred to as *Evidence Theory*, is based on representing independent pieces of evidence by mass functions and combining them using a generic operator called Dempster’s rule. It is a well-established formalism for reasoning and making decisions under uncertainty [10, 13, 38]. A mass function has more degrees of freedom than a probability distribution, which allows it to represent a wider range of belief states, from complete ignorance to full certainty.

One of the applications of the DS theory is to design *evidential classifiers* (e.g., [9, 14, 24, 33]), which compute a predicted mass function for each input vector. The output mass function can then be used for decision-making [4, 8, 19]. Over the years, two main principles have been developed to design an evidential classifier: the *model-based* and *distance-based* approaches [15]. The former uses estimated class-conditional distributions [32], while the latter constructs mass functions based on distances to prototypes [9, 14, 24]. Thanks to the generality and

expressiveness of the belief-function formalism, an evidential classifier provides more informative outputs than those of conventional classifiers (e.g., a neural network with a softmax output layer). The flexibility of evidential classifiers can be exploited for uncertain data classification [40] and set-valued classification [8, 26]. Therefore, it may be advantageous to combine an FCN-based model with an evidential classifier for semantic segmentation.

The objective of this study is to take advantage of object representations generated by an FCN and use them as the input features of an evidential classifier for decision-making. The proposed model, referred to as the *evidential fully convolutional network* (E-FCN), transforms an FCN model by replacing its softmax layer by a distance-based DS layer and a utility layer. In an E-FCN, an FCN architecture is used to extract pixel-wise high-order features from an input image. Then, the features are converted into pixel-wise mass functions by the DS layer. Finally, the mass functions are used to compute the utilities of acts assigning pixels to a set of classes for semantic segmentation in a so-called “utility layer”. An end-to-end learning procedure allows us to train the E-FCN using a learning set with soft labels. The effectiveness of the E-FCN is demonstrated and discussed based on experiments using three benchmark databases (Pascal VOC 2011 [17], MIT-scene Parsing [43], and SIFT Flow [34]).

The rest of the paper is organized as follows. Section 2 starts with a brief reminder of DS theory, the DS layer for constructing mass functions, and feature representation via FCN. The E-FCN model is then introduced in Section 3. Section 4 presents numerical experiments, which demonstrate the advantages of the E-FCNs. Finally, we conclude the paper in Section 5.

2 Background

This section first recalls some necessary definitions regarding DS theory (Section 2.1), and the evidential neural network at the basis of the DS layer (Section 2.2). A brief description of feature representation via FCNs is then provided in Section 2.3.

2.1 Dempster-Shafer theory

The main concepts underlying DS theory are only briefly presented in this section, and some basic notations are introduced. Detailed information can be found in Shafer’s original work [30] and in the recent review [13].

Let $\Omega = \{\omega_1, \dots, \omega_M\}$ be a set of classes, called the *frame of discernment*. A *mass function* on Ω is a mapping m from 2^Ω to $[0,1]$ such that $m(\emptyset) = 0$ and

$$\sum_{A \subseteq \Omega} m(A) = 1. \quad (1)$$

For any $A \subseteq \Omega$, each mass $m(A)$ is interpreted as a share of a unit mass of belief allocated to the hypothesis that the truth is in A , and which cannot be allocated to any strict subset of A based on the available evidence. Set A is called a *focal set* of m if $m(A) > 0$. A mass function is said to be *logical* if it has only one focal set.

Two mass functions m_1 and m_2 representing independent items of evidence can be combined conjunctively by Dempster's rule \oplus [30] as

$$(m_1 \oplus m_2)(A) = \frac{(m_1 \cap m_2)(A)}{1 - (m_1 \cap m_2)(\emptyset)} \quad (2a)$$

for all $A \neq \emptyset$, with

$$(m_1 \cap m_2)(A) = \sum_{B \cap C = A} m_1(B) m_2(C), \quad (2b)$$

and $(m_1 \oplus m_2)(\emptyset) = 0$. Mass functions m_1 and m_2 can be combined if and only if the denominator on the right-hand side of Eq. (2a) is strictly positive. The operator \oplus is commutative and associative.

For decision-making with belief functions, let $u_{ij} \in [0, 1]$ denote the utility of selecting ω_i when the true state is ω_j , and f_{ω_i} the act of selecting ω_i . We define the *pignistic expected utility* [11] of act f_{ω_i} as

$$\mathbb{E}_m(f_{\omega_i}) = \sum_{j=1}^M u_{ij} \text{Bet}P_m(\{\omega_j\}), \quad (3a)$$

where $\text{Bet}P_m$ is the *pignistic* probability measure computed from mass function m by the *pignistic transformation*, defined as

$$\text{Bet}P_m(\{\omega_j\}) = \sum_{\{A \subseteq \Omega: \omega_j \in A\}} \frac{m(A)}{|A|}, \quad (3b)$$

for all $\omega_j \in \Omega$. The act with the highest pignistic expected utility can then be selected. Other decision criteria in the belief function framework are reviewed in [11] and [26].

2.2 Evidential neural network

Denœux [9] proposed a distance-based neural-network based on DS theory, known as the *evidential neural network (ENN) classifier*. The ENN classifier summarizes the learning set by a small number of prototypes, and treats the proximity of an input vector to each prototype as a piece of evidence about its class. The different pieces of evidence are represented by mass functions, which are combined using Dempster's rule (2). This section provides a brief description of the ENN classifier.

We consider a training set $\mathcal{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\} \subset \mathbb{R}^P$ of N examples represented by P -dimensional feature vectors, and n prototypes $\{\mathbf{p}^1, \dots, \mathbf{p}^n\} \subset \mathbb{R}^P$. For a test sample \mathbf{x} , the ENN classifier constructs mass functions that quantify the uncertainty about its class in $\Omega = \{\omega_1, \dots, \omega_M\}$, using a three-step procedure. This procedure can be implemented in a neural network with two hidden layers and one output layer. These three layers can be considered as a single complex layer called the "DS layer", which will be plugged into an FCN architecture as explained in Section 3.1. The three-step procedure can be described as follows.

125 Step 1: The similarity between \mathbf{x} and each prototype \mathbf{p}^l is computed as

$$s^l = \alpha^l \exp\left(-\left(\eta^l d^l\right)^2\right), \quad l = 1, \dots, n, \quad (4)$$

126 where $d^l = \|\mathbf{x} - \mathbf{p}^l\|$ is the Euclidean distance between \mathbf{x} and prototype \mathbf{p}^l ,
 127 $\eta^l \in \mathbb{R}$ is a scale parameter and α^l is a parameter in $[0, 1]$. Prototypes $\mathbf{p}^1, \dots, \mathbf{p}^n$
 128 can be considered as vectors of connection weights between the input layer
 129 and a hidden layer of n Radial Basis Function (RBF) units. **The number n of**
 130 **prototypes is a hyper-parameter and can be tuned using a validation set or by**
 131 **cross-validation.**

Step 2: The mass function m^l associated to reference pattern \mathbf{p}^l is computed as

$$m^l(\{\omega_j\}) = v_j^l s^l, \quad j = 1, \dots, M \quad (5a)$$

$$m^l(\Omega) = 1 - s^l, \quad (5b)$$

132 where $v_j^l \geq 0$ is the degree of membership of prototype \mathbf{p}^l to class ω_j with
 133 $\sum_{j=1}^M v_j^l = 1$. We denote the vector of masses induced by prototype \mathbf{p}^l as

$$\mathbf{m}^l = (m^l(\{\omega_1\}), \dots, m^l(\{\omega_M\}), m^l(\Omega))^T.$$

134 Eq. (5) can be regarded as computing the activation of units in a “mass func-
 135 tions” layer composed of n modules of $M + 1$ units each. The activations of
 136 the units in module l correspond to the belief masses assigned by m^l .

Step 3: The n mass functions \mathbf{m}^l , $l = 1, \dots, n$, are aggregated by Dempster’s rule
 (2). The combined mass function can be computed iteratively as $\mu^1 = \mathbf{m}^1$ and
 $\mu^l = \mu^{l-1} \cap m^l$ for $l = 2, \dots, n$. From (2a), we have

$$\mu^l(\{\omega_j\}) = \mu^{l-1}(\{\omega_j\})m^l(\{\omega_j\}) + \mu^{l-1}(\{\omega_j\})m^l(\{\Omega\}) + \mu^{l-1}(\Omega)m^l(\{\omega_j\}) \quad (6a)$$

137 for $l = 2, \dots, n$ and $j = 1, \dots, M$, and

$$\mu^l(\Omega) = \mu^{l-1}(\Omega)m^l(\Omega) \quad l = 2, \dots, n. \quad (6b)$$

138 The output vector $\mathbf{m} = (m(\{\omega_1\}), \dots, m(\{\omega_M\}), m(\Omega))^T$ is finally obtained
 139 by normalizing μ^n as

$$m(A) = \frac{\mu^n(A)}{\mu^n(\Omega) + \sum_{j'=1}^M \mu^n(\{\omega_{j'}\})},$$

140 with $A \in \{\{\omega_1\}, \dots, \{\omega_M\}, \Omega\}$.

141 The network parameters are the prototypes \mathbf{p}^l , the coefficients α^l and η_l , and
 142 the membership degrees v_j^l for $l = 1, \dots, n$ and $j = 1, \dots, M$. They are learnt by
 143 minimizing a loss function. To enforce the constraints $0 \leq \alpha^l \leq 1$, we introduce
 144 new variables $\xi^l \in \mathbb{R}$ such that

$$\alpha^l = \frac{1}{1 + \exp(-\xi^l)} \in (0, 1).$$

145 Similarly, the constraints on parameters v_j^l are enforced by introducing new pa-
 146 rameters $\delta_j^l \in \mathbb{R}$ such that

$$v_j^l = \frac{(\delta_j^l)^2}{\sum_{j'=1}^M (\delta_{j'}^l)^2} \quad (7)$$

147 for $l = 1, \dots, n$ and $j = 1, \dots, M$. More details can be found in [9].

148 2.3 Fully convolutional network

149 The performance of a classifier in semantic segmentation tasks heavily depends on
 150 the information contained in its input features. Feature representation, an essential
 151 part of the machine learning workflow, consists in discovering the predictors needed
 152 for semantic segmentation from input images. In recent years, FCNs [25] and their
 153 variants [22, 28] have achieved remarkable performances thanks to their ability to
 154 construct rich pixel-wise deep feature representations.

155 FCNs owe their name to their architecture, which is built only from locally
 156 connected layers, such as convolution, pooling, and upsampling layers. No dense
 157 layer is used in this kind of architecture. Generally, an FCN consists of two main
 158 parts: an encoder-decoder architecture for pixel-wise object representation and
 159 a softmax layer for pixel-wise assignments. In the encoder-decoder architecture,
 160 an input image is encoded by several convolutional and pooling layers and then
 161 decoded by one or more upsampling layers. The softmax layer assigns each pixel in
 162 the input image to one of the classes based on the outputs of the encoder-decoder
 163 architecture. Therefore, the outputs of the encoder-decoder architecture, called
 164 the *pixel-wise feature maps*, are considered as a feature representation of the input
 165 image. In the study, these feature maps are used as input to a DS layer allowing
 166 for set-valued semantic segmentation, as will be shown in Section 3.1.

167 To understand the feature representation of FCNs, we briefly recall the encoder-
 168 decoder architecture illustrated in Figure 1. The encoder part consists of several
 169 convolutional and pooling layers. Each convolutional layer performs convolutions
 170 its input to produce a set of feature maps. Let $\mathbf{z} = (z^1, \dots, z^D)$ be the input
 171 made up of D input maps or *input channels* z^i ($i = 1, \dots, D$) of size $H \times W$. The
 172 processes in a convolutional layer with input \mathbf{z} , consisting of e convolution kernels
 173 with size $a \times b$, are expressed as

$$c^j = f(\lambda^j + \sum_i \varepsilon^{i,j} * z^i), \quad (8)$$

174 where $\varepsilon^{i,j}$, a matrix of size $a \times b$, is the convolution kernel between the i -th input
 175 map and the j -th output map; λ^j is the bias of kernel $\varepsilon^{i,j}$; $*$ denotes the convolution
 176 operation; c^j is the j -th output feature map, with size $\frac{h-a+1}{r} \times \frac{w-b+1}{r}$, $j = 1, \dots, e$;
 177 r is the stride with which the kernel slides over input map z^i , and f is the activation
 178 function, such as the rectified linear unit $\text{ReLU}(x) = \max(0, x)$ [23]. A pooling layer
 179 follows the convolutional layer to sub-sample feature map c^j by computing some
 180 statistics of feature values within non-overlapping $s \times s$ windows. In the case of
 181 max-pooling used in this paper, the statistic is the maximum. Thus, the outputs of
 182 the pooling layer is composed of the D feature maps sub-sampled by factor s . For
 183 example, feature map c^j with size $\frac{h-a+1}{r} \times \frac{w-b+1}{r}$ is downsized to $\frac{h-a+1}{2r} \times \frac{w-b+1}{2r}$
 184 by a pooling layer with a 2×2 non-overlapping window.

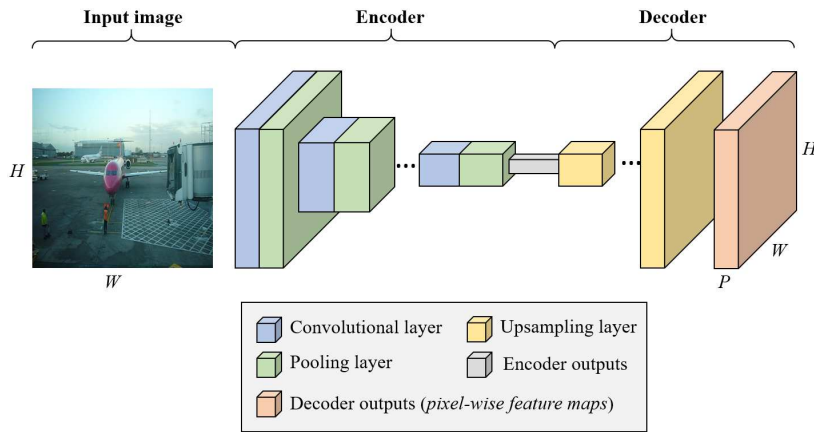


Fig. 1: Illustration of the encoder-decoder architecture. An encoder downsizes its input by convolution and pooling operations. The outputs of the encoder, as the sparse feature maps, are imported into a decoder. A decoder upsamples and densifies its inputs by performing the reverse operation of convolution and pooling. The final decoder outputs are the pixel-wise feature maps.

185 Although the convolution and pooling operations in the encoder part contribute
 186 to feature representation by retaining only robust activations, spatial information
 187 within a receptive field is lost, which may be critical for image semantic segmen-
 188 tation. To address the issue, a decoder part made up of one or more upsampling
 189 layers is added at the output of the encoder part. The decoder performs the re-
 190 verse operation of convolution and pooling for reconstructing a set of activations
 191 with the same size of the input image, as shown in Figure 1. Thus, the outputs of
 192 the decoder part are enlarged feature maps. In this study, we use a *deconvolution*
 193 *layer* [28] to implement the upsampling operation.

194 A deconvolutional layer densifies its inputs of sparse feature maps through
 195 convolution-like operations with multiple learned kernels. However, contrary to
 196 convolutional layers, which connect multiple inputs within a kernel to a single
 197 activation, a deconvolutional layer associates a single input in a feature map to
 198 multiple outputs. Thus, the outputs of a deconvolutional layer are enlarged and
 199 dense feature maps. The processes of a deconvolution operation can also be sum-
 200 marized as Eq. (8), but its kernel sizes are larger than the input sizes, i.e., $a \geq H$
 201 and $b \geq W$.

202 3 Evidential fully convolutional network

203 In this section, we describe the proposed E-FCN. Section 3.1 presents the overall
 204 architecture composed of an encoder-decoder module for feature representation, a
 205 DS layer to construct mass functions, and a utility layer for decision-making. The
 206 details of the utility layer are described in Section 3.2. Section 3.3 introduces the
 207 strategy for training E-FCN models using a learning set with soft labels.

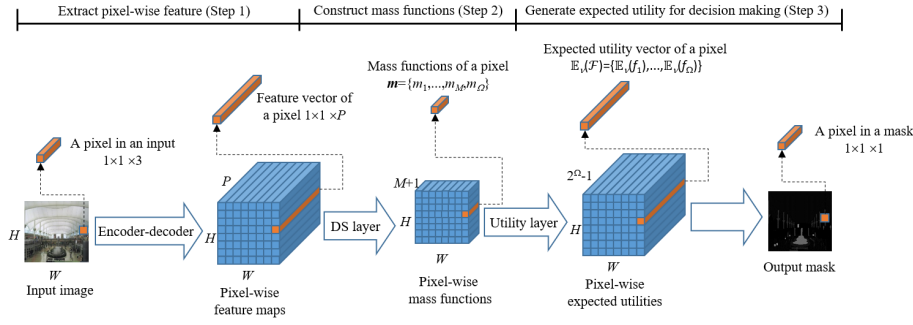


Fig. 2: Architecture of an evidential fully convolutional network (E-FCN). The E-FCN performs semantic segmentation using a three-step procedure. In the first step, an encoder-decoder architecture extracts pixel-wise feature maps from the input image. Each vector in the feature maps is fed into a DS layer to construct the pixel-wise mass functions in the second step. These mass functions are finally fed into a utility layer to generate the pixel-wise expected utilities of all acts. Finally, the segmentation mask is computed based on the expected utilities.

208 3.1 Network architecture

209 The main idea of this work is to hybridize the ENN classifier presented in Section 2.2 and the FCN recalled in Section 2.3 by “plugging” a DS layer followed
 210 by a utility layer at the output of the final deconvolutional layer in the FCN.
 211 The architecture of the proposed method, called the *evidential FCN* (E-FCN), is
 212 illustrated in Figure 2. An E-FCN classifier performs set-valued semantic segmen-
 213 tation and quantifies the uncertainty about the class of each pixel, taking values
 214 in $\Omega = \{\omega_1, \dots, \omega_M\}$, using a three-step procedure defined as follows.
 215

216 Step 1: As in a probabilistic FCN (P-FCN), an image of size $W \times H \times 3$ is
 217 presented as input to the encoder-decoder architecture of an FCN to generate
 218 pixel-wise feature maps of size $W \times H \times P$, where P is the number of *output*
 219 *channels*. Each feature vector $1 \times 1 \times P$ from a pixel-wise feature map is a
 220 P -dimensional representation of the corresponding pixel, ready to be fed into
 221 the DS layer. This architecture generates reliable pixel-wise representations of
 222 the input image. Thanks to the representations, the E-FCN yields similar or
 223 even better performance for precise semantic segmentation than does a P-FCN
 224 with the same encoder-decoder architecture, as will be shown in Section 4.2.

225 Step 2: Each feature vector from the encoder-decoder architecture is fed into
 226 the DS layer, in which it is used to compute a mass function as explained in
 227 Section 2.2. The output of the DS layer for a given feature vector is an
 228 $(M + 1)$ -dimensional mass vector

$$\mathbf{m} = (m(\{\omega_1\}), \dots, m(\{\omega_M\}), m(\Omega))^T.$$

229 Thus, given pixel-wise feature maps of size $W \times H \times P$ from Step 1, the output
 230 of the DS layer is a tensor of size $W \times H \times (M + 1)$. Each mass vector in
 231 the tensor represents the uncertainty about the class of the corresponding
 232 pixel. More precisely, the mass $m(\{\omega_i\})$ is the degree of belief that the true

class of the corresponding pixel is ω_i . The DS layer tends to allocate uniform masses if the representations contain confusing information. The additional degree of freedom $m(\Omega)$ makes it possible to quantify the lack of evidence [12] and to verify whether the model is well trained [35]. The advantages of this uncertainty representation will be demonstrated in the performance evaluation of set-valued semantic segmentation using E-FCN in Section 4.3.

Step 3: The output pixel-wise mass vectors are fed into a utility layer, where they are used to compute the expected utility of acts. Each act is defined as the assignment of a pixel to a non-empty subset A of Ω . Therefore, the output of the layer for each feature vector from Step 2 is a vector of at most $2^\Omega - 1$ expected utilities when all of the possible acts are considered. The utility layer allows the E-FCN to perform set-valued semantic segmentation. This capability will be demonstrated by comparing the performances of the two types of FCNs in set-valued segmentation (Section 4.3) and novelty detection (Section 4.4) tasks. More details about the utility layer are given in the next section.

3.2 Utility layer

In this section, we describe in greater detail the decision-making process taking place in the utility layer. Section 3.2 begins by introducing the precise semantic segmentation method using mass functions and utilities. Section 3.2 then describes a method for computing the utility of set-valued pixel-wise classification, after which an approach to set-valued classification based on mass functions is described in Section 3.2. In Section 3.2, we summarize the workflow as a neural network layer in the E-FCN model.

3.2.1 Precise semantic segmentation

Let $\Omega = \{\omega_1, \dots, \omega_M\}$ be the set of classes. For semantic segmentation problems with precise prediction, each pixel in an image is assigned to exactly one class. An act is thus defined as the assignment of a pixel to one and only one of the M classes, and the set of acts is $\mathcal{F} = \{f_{\omega_1}, \dots, f_{\omega_M}\}$, where f_{ω_i} denotes assignment to class ω_i . To make decisions, we define a utility matrix \mathbf{U} of size $M \times M$, whose general term $u_{ij} \in [0, 1]$ is the utility of assigning a pixel to class ω_i when the true class is ω_j .

When uncertainty about Ω is described by a DS mass function, each act f_{ω_i} induces an expected utility, such as the pignistic expected utility defined by Eq. (3). Given utility matrix \mathbf{U} and the output of the DS layer \mathbf{m} for a given pixel, the pignistic expected utility of assigning that pixel to class ω_i is given by Eq. (3a), where $BetP_m$ is the pignistic probability defined by Eq. (3b). The pixel is finally assigned to class ω_{i^*} such that

$$i^* = \arg \max_{i \in \{1, \dots, M\}} \mathbb{E}_m(f_{\omega_i}). \quad (9)$$

3.2.2 Extending the utility matrix

For semantic segmentation problems with imprecise prediction, we adopt the approach described in [26] for set-valued classification under uncertainty, which allows the assignment of a pixel to any non-empty subset A of Ω . The set of acts thus potentially becomes $\mathcal{F} = \{f_A : A \subseteq \Omega, A \neq \emptyset\}$, where f_A denotes the assignment to a subset A . (In practice, when the cardinality of Ω is very large, we may only consider acts f_A for *some* subsets A of Ω). In this study, f_A is referred to as an *imprecise assignment* when subset A is a *multi-class set* with $|A| \geq 2$. For decision-making with \mathcal{F} , the utility matrix \mathbf{U} has to be extended to a matrix $\tilde{\mathbf{U}}$ of size $(2^M - 1) \times M$, where each element $\tilde{u}_{A,j}$ denotes the utility of assigning a pixel to set A of classes when the true label is ω_j . Following [26], this extension is performed as follows.

When the true class is ω_j , the utility of assigning a pixel to set A is defined as an *Ordered Weighted Average (OWA)* aggregation [37] of the utilities of each precise assignment in A as

$$\tilde{u}_{A,j} = \sum_{k=1}^{|A|} g_k u_{(k)j}^A, \quad (10)$$

where $u_{(k)j}^A$ is the k -th largest element in the set $\{u_{ij} : \omega_i \in A\}$ made up of the elements in the utility matrix \mathbf{U} , and weights $\mathbf{g} = (g_1, \dots, g_{|A|})$ represent the preference to choose $u_{(k)j}^A$ if forced to select a single value in $\{u_{ij} : \omega_i \in A\}$. The components of weight vector \mathbf{g} represent the *tolerance to imprecision* of a decision maker (DM). For example, full tolerance to imprecision is achieved when the assignment act f_A has utility 1 once set A contains the true label, no matter how large A is. In this case, only the maximum utility of elements in set $\{u_{ij}, \omega_i \in A\}$ is considered: $(g_1, g_2, \dots, g_{|A|}) = (1, 0, \dots, 0)$. At the other extreme, a DM attaching no value to imprecision would consider the act f_A as equivalent to selecting one class uniformly at random from A : this is achieved when

$$(g_1, g_2, \dots, g_{|A|}) = \left(\frac{1}{|A|}, \frac{1}{|A|}, \dots, \frac{1}{|A|} \right),$$

in which case the OWA operator becomes the average. In this study, following [26], we determine the weight vector \mathbf{g} of the OWA operator by adapting O'Hagan's method [29]. We define the tolerance to imprecision as

$$TDI(\mathbf{g}) = \sum_{k=1}^{|A|} \frac{|A| - k}{|A| - 1} g_k = \gamma, \quad (11)$$

which equals 1 for the maximum, 0 for the minimum, and 0.5 for the average. In practice, we only need to consider values of γ between 0.5 and 1 as a precise assignment is always more desirable than an imprecise one when $\gamma < 0.5$ [26]. Given a value of γ , we can compute the weights of the OWA operator by maximizing the entropy

$$ENT(\mathbf{g}) = - \sum_{k=1}^{|A|} g_k \log g_k \quad (12)$$

subject to the constraints $TDI(\mathbf{g}) = \gamma$, $\sum_{k=1}^{|A|} g_k = 1$, and $g_k \geq 0$.

Table 1: Utility matrix extended by an OWA operator with $\gamma = 0.8$.

	Classes		
	ω_1	ω_2	ω_3
$f_{\{\omega_1\}}$	1	0	0
$f_{\{\omega_2\}}$	0	1	0
$f_{\{\omega_3\}}$	0	0	1
$f_{\{\omega_1, \omega_2\}}$	0.8	0.8	0
$f_{\{\omega_1, \omega_3\}}$	0.8	0	0.8
$f_{\{\omega_2, \omega_3\}}$	0	0.8	0.8
f_Ω	0.6819	0.6819	0.6819

305 **Example 1** Table 1 shows an example of the extended utility matrix generated by
306 an OWA operator with $\gamma = 0.8$. The first three rows constitute the original utility
307 matrix, indicating that the utility equals 1 when assigning a sample to its true class,
308 and 0 otherwise. The remaining rows are the matrix of the aggregated utilities. For
309 example, we get a utility of 0.8 when assigning a sample to set $\{\omega_1, \omega_2\}$ if the true
310 label is ω_1 .

311 3.2.3 Set-valued semantic segmentation

312 Based on an extended utility matrix \tilde{U} and the output of the DS layer \mathbf{m} for
313 a given pixel, we can compute the pignistic expected utility of assigning that pixel
314 to set A as

$$\mathbb{E}_m(f_A) = \sum_{j=1}^M \tilde{u}_{A,j} \text{Bet}P_m(\{\omega_j\}), \quad (13)$$

315 where $\text{Bet}P_m$ is the pignistic probability defined by Eq. (3b). The pixel is finally
316 assigned to set A such that

$$A = \arg \max_{\emptyset \neq B \subseteq \Omega} \mathbb{E}_m(f_B). \quad (14)$$

317 3.2.4 Utility layer

318 The procedure of assigning a pixel to a set of classes using utility theory is
319 implemented in the neural network as a *utility layer*. In this layer, the inputs
320 and outputs are, respectively, the pixel-wise mass vectors \mathbf{m} from the preceding
321 DS layer and the pixel-wise expected utilities of all acts in \mathcal{F} . The connection
322 weight between unit j of the DS layer and output unit $A \subseteq \Omega$ corresponding
323 to the assignment to set A is the utility value $\tilde{u}_{A,j}$. As coefficient γ describing
324 the imprecision tolerance degree is predetermined, the connection weights of the
325 expected utility layer are fixed and do not need to be updated during training.

326 In practice, the connections between the DS and utility layers can be designed
327 by the user. For example, one can build a utility layer using the utility values $\tilde{u}_{A,j}$
328 with $|A| = 1$ to only consider precise assignments, or $0 < |A| \leq 2$ to consider
329 assignment to sets classes of cardinality one or two. In this paper, we have only

330 considered the acts f_A such that A is a singleton, Ω , or one of the soft labels
 331 present in the learning set (as explained in Section 3.3 below).

332 3.3 Learning with soft labels

333 In traditional learning systems for image semantic segmentation, all pixels are la-
 334 beled with a single class even when their true class cannot be determined with full
 335 certainty. For example, the true class may be uncertain at object borders, but the
 336 border pixels are still given precise labels. Additionally, one cannot reliably label
 337 some small objects in an image, such as distant objects in a driving scene. Arbi-
 338 trarily giving precise labels to pixels with confusing information may negatively
 339 impact the performance of learning systems in image semantic segmentation tasks.
 340 The notion of *soft label* [5, 14] is a way to address this problem.

341 Here, we define a soft label as a nonempty subset $A_* \in 2^\Omega \setminus \emptyset$ of classes a pixel
 342 may belong to, based on our current knowledge. For example, label $A_* = \{\omega_i, \omega_j\}$
 343 indicates that the true class of a pixel is known to be either ω_i or ω_j but we cannot
 344 determine which one specifically. A strategy of end-to-end learning is proposed to
 345 train an E-FNC from an image learning set with soft labels. All parameters in the
 346 DS layer are first initialized randomly using normal distributions. For a given pixel
 347 with nonempty soft label $A_* \subseteq \Omega$, let m_l be the logical mass function with focal set
 348 A_* , i.e., such that $m_l(A_*) = 1$. The *labeling* pignistic expected utilities $\mathbb{E}_{m_l}(f_A)$
 349 for $A \in 2^\Omega \setminus \emptyset$ can be computed using Eq. (13) and the pignistic belief-probability
 350 transformation Eq. (3b). Similarly, we consider the *predicted* pignistic expected
 351 utilities $\mathbb{E}_m(f_A)$ for $A \in 2^\Omega \setminus \emptyset$, where m is the predicted mass function from the
 352 DS layer of the E-FNC, with focal sets $\{\omega_1\}, \dots, \{\omega_M\}$ and Ω . For a given pixel
 353 with soft label m_l and predicted mass function m , the loss $\mathcal{L}(m, m_l)$ is defined as
 354 the squared Euclidean distance between the vectors of expected utilities w.r.t. m_l
 355 and m :

$$\mathcal{L}(m, m_l) = \sum_{\emptyset \neq A \subseteq \Omega} [\mathbb{E}_{m_l}(f_A) - \mathbb{E}_m(f_A)]^2. \quad (15)$$

356 The derivatives of $\mathcal{L}(m, m_l)$ of the error w.r.t the output masses $m(\{\omega_k\})$ are

$$\begin{aligned} \frac{\partial \mathcal{L}(m, m_l)}{\partial m(\{\omega_k\})} &= \sum_{\emptyset \neq A \subseteq \Omega} \frac{\partial \mathcal{L}(m, m_l)}{\partial \mathbb{E}_m(f_A)} \cdot \frac{\partial \mathbb{E}_m(f_A)}{\partial m(\{\omega_k\})} \\ &= -2 \sum_{\emptyset \neq A \subseteq \Omega} [\mathbb{E}_{m_l}(f_A) - \mathbb{E}_m(f_A)] \sum_{j=1}^M \frac{\partial \mathbb{E}_m(f_A)}{\partial \text{Bet}P_m(\omega_j)} \frac{\partial \text{Bet}P_m(\omega_j)}{\partial m(\{\omega_k\})} \\ &= -2 \sum_{\emptyset \neq A \subseteq \Omega} [\mathbb{E}_{m_l}(f_A) - \mathbb{E}_m(f_A)] \sum_{j=1}^M \tilde{u}_{A,j} \left(\delta_{kj} - \frac{1}{M} \right), \end{aligned} \quad (16)$$

357 where $\delta_{kj} = 1$ if $k = j$ and $\delta_{kj} = 0$ otherwise. The derivatives of $m(\{\omega_k\})$ w.r.t
 358 p_k^l , η^l , and ξ^l in the DS layer are the same as in Denœux's original work [9], and
 359 the gradient with respect to all network parameters can be back-propagated from
 360 the output layer to the input layer.

4 Experiments

In this section, we present numerical experiments that demonstrate the advantages of the proposed model. The databases and metrics are first introduced in Section 4.1. Precise and imprecise segmentation results are then reported, respectively, in Sections 4.2 and 4.3. Finally, novelty detection results are presented in Section 4.4.

4.1 Databases and metrics for performance evaluation

Databases

Three benchmark databases were used in the study: Pascal VOC 2011 [17], MIT-scene Parsing [43], and SIFT Flow [34]. These databases were used to train and test the E-FCNs as well as probabilistic FCNs (P-FCNs) for comparison.

The Pascal VOC 2011 database contains 20 object classes in 5034 images, with segmentation masks that indicate the class of each pixel, or label it as “background” if the object does not belong to one of the twenty specified classes. The MIT-scene Parsing and SIFT Flow databases are similar to the Pascal VOC 2011 database but have, respectively, 150 categories in 20K images and 33 classes in 2688 images. The list of classes for the three databases are given in Table 2. Each of the three databases was split into 50% for training/validation and 50% for testing. The validation sets were used to determine hyper-parameters, such as the number of prototypes in each DS layer. The optimal tolerance to imprecision γ can be determined in the same way since it can also be considered as a hyper-parameter.

There is no confidence value associated with the pixel labels in any of the three databases. Thus, we defined soft labels for them. For the Pascal VOC 2011 database, we assigned each pixel in a boundary area a soft label $A \subseteq \Omega$, where A consists of the object classes around the boundary area. Some examples are shown in Figure 3a. For the MIT-scene Parsing and SIFT Flow databases with no identified boundary areas, we assigned soft labels to the pixels situated between every two objects, as shown in Figures 3b and 3c.

A semantic segmentation model should not only be accurate for the classes in the learning set, but it should also be able to detect objects belonging to classes that are not included in the learning set. To evaluate this novelty detection capacity, we mixed the three databases: for example, an FCN model trained using the Pascal VOC 2011 database was tested on the other two databases.

Metrics

We used three metrics for the performance evaluation of semantic segmentation: pixel utility (PU), utility of intersection over union (UIoU), and expected calibration error (ECE).

Pixel utility. For an image with T pixels, the *pixel utility* is defined as

$$PU = \frac{1}{|T|} \sum_{i=1}^{|T|} \tilde{u}_{A(i), A_*(i)} \quad (17)$$

Table 2: Lists of classes for the three databases used in this study. Classes in bold characters are included in two or three databases. Classes with close meanings, such as “minibike” and “motorbike”, are considered as identical.

Database	Class list
Pascal VOC 2011	background, cat, dog, horse, sheep, train, sofa, aeroplane, bicycle, bird, boat, bottle, bus, car, chair, cow, diningtable, motorbike, person, pottedplant, tv.
MIT-scene parsing	wall, floor, ceiling, bed, cabinet, earth, curtain, water, painting, shelf, house, mirror, rug, armchair, seat, desk, wardrobe, lamp, bathtub, railing, cushion, base, box, column, chest, counter, sink, skyscraper, fireplace, refrigerator, grandstand, path, stairs, runway, case, pool, pillow, screen, bookcase, blind, coffee, toilet, flower, book, hill, bench, countertop, stove, palm, kitchen, computer, swivel, bar, arcade, hovel, towel, light, truck, tower, chandelier, booth, dirt track, apparel, land, banner, escalator, ottoman, buffet, poster, stage, van, ship, fountain, conveyer, canopy, washer, plaything, swimming, stool, barrel, basket, waterfall, tent, bag, minibike, cradle, oven, ball, food, step, tank, trade, microwave, pot, animal, lake, dishwasher, screen, blanket, sculpture, hood, sconce, vase, traffic, tray, ashcan, fan, pier, screen, plate, monitor, bulletin, shower, radiator, glass, clock, flag, sofa, airplane, building, sky, tree, road, windowpane, grass, sidewalk, person, door, table, mountain, plant, chair, car, sea, field, fence, rock, sign, sand, staircase, river, bridge, boat, bus, awning, streetlight, tv, pole, bottle, minibike, bicycle.
SIFT Flow	balcony, crosswalk, desert, moon, sun, window, awning, bird, boat, bridge, building, bus, car, cow, door, fence, field, grass, mountain, person, plant, pole, river, road, rock, sand, sea, sidewalk, sign, sky, staircase, streetlight, tree.

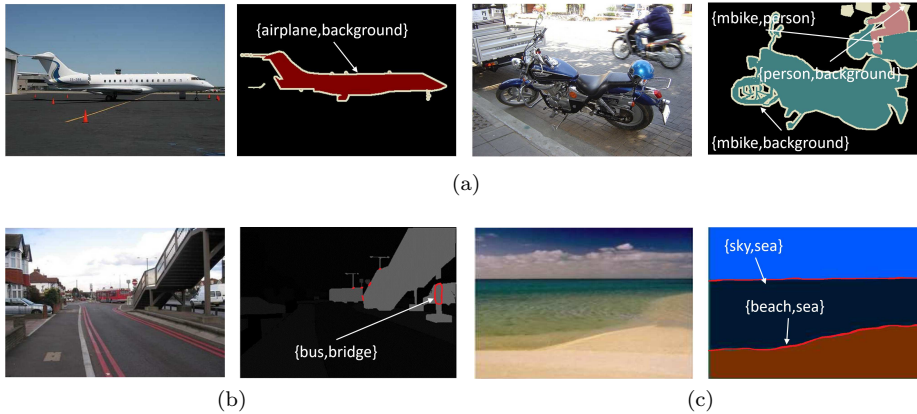


Fig. 3: Segmentation masks with soft labels: (a) Pascal VOC 2011, (b) MIT-scene Parsing, and (c) SIFT Flow.

398 where $A_*(i)$ is the label of pixel i , $A(i)$ is the selected set of classes for pixel
 399 i determined from Eq. (14), and using the notations introduced in Section 3.2,
 400 $\tilde{u}_{A(i), A_*(i)}$ is the utility of assigning pixel i to subset $A(i) \subseteq \Omega$ when its label
 401 is $A_*(i)$. Thus, PU is the same as pixel accuracy when only considering precise
 402 assignments and precise labels. To consider soft labels, the utility matrix \tilde{U} defined

Table 3: Utility matrix considering soft labels with $\gamma = 0.8$.

		Label						
		ω_1	ω_2	ω_3	$\{\omega_1, \omega_2\}$	$\{\omega_1, \omega_3\}$	$\{\omega_2, \omega_3\}$	Ω
Act	$f_{\{\omega_1\}}$	1	0	0	0.625	0.625	0	0.489
	$f_{\{\omega_2\}}$	0	1	0	0.625	0	0.625	0.489
	$f_{\{\omega_3\}}$	0	0	1	0	0.625	0.625	0.489
	$f_{\{\omega_1, \omega_2\}}$	0.8	0.8	0	1	0.5	0.5	0.782
	$f_{\{\omega_1, \omega_3\}}$	0.8	0	0.8	0.5	1	0.5	0.782
	$f_{\{\omega_2, \omega_3\}}$	0	0.8	0.8	0.5	0.5	1	0.782
	f_{Ω}	0.682	0.682	0.682	0.853	0.853	0.853	1

403 in Section 3.2 should be extended to a matrix $\tilde{\mathbf{U}}'$ of size $(2^M - 1) \times (2^M - 1)$ with
 404 general term \tilde{u}_{A, A_*} defined as the utility of assigning a pixel to subset $A \subseteq \Omega$
 405 when its label is A_* , with $|A_*| \geq 1$. Soft label A_* means that we only know the
 406 true class of a pixel is in set A_* , and nothing more. To define the utility \tilde{u}_{A, A_*} , we
 407 first compute the average of the utilities of selecting subset A when the true class
 408 is in A_* as

$$\bar{u}_{A, A_*} = \frac{1}{|A_*|} \sum_{\omega_k \in A_*} \tilde{u}_{A, k}, \quad (18a)$$

409 where $\tilde{u}_{A, k}$ is the utility of selecting subset A when the true class is ω_k , and we
 410 normalize this average utility to ensure that $\tilde{u}_{A_*, A_*} = 1$:

$$\tilde{u}_{A, A_*} = \frac{\bar{u}_{A, A_*}}{\bar{u}_{A_*, A_*}}. \quad (18b)$$

411

412 **Example 2** Table 3 shows an example of the utility matrix considering soft labels,
 413 which is extended from Example 1. The last four columns correspond to the utility
 414 matrix for soft labels. An act achieves utility 1 only if $A = A_*$, 0 if $A \cap A_* = \emptyset$,
 415 and a value between 0 and 1 if $A \neq A_*$ and $A \cap A_* \neq \emptyset$.

416 *Utility of intersection over union.* The segmentation performance was also evalu-
 417 ated by the utility of intersection over union (UIoU) defined as

$$UIoU = \frac{1}{2^{|\Omega|} - 1} \sum_{B \subseteq \Omega} \frac{\sum_{i \in \mathbf{G}^B \cap \mathbf{P}^B} \tilde{u}_{A(i), B}}{|\mathbf{G}^B \cup \mathbf{P}^B|}, \quad (19)$$

418 where $\mathbf{P}^B = \{i : A(i) \cap B \neq \emptyset\}$ is the predicted area containing pixels assigned
 419 to a set of classes that intersect B , and $\mathbf{G}^B = \{i : A_*(i) = B\}$ is the ground
 420 truth area composed of pixels with label B . Thus, in the special case of precise
 421 segmentation with only precise labels, UIoU boils down to *intersection over union*,
 422 a widely used metric for semantic segmentation [22, 25, 28].

423 *Expected calibration error.* In decision systems, a neural network should not only
 424 be accurate, but it should also indicate when it is likely to be incorrect. Thus,
 425 the confidence of an E-FCN should be *calibrated*. To characterize this property,

426 we extend the *expected calibration error* (ECE) introduced in [20] as follows. We
 427 define the *prediction confidence* of pixel i as

$$co(i) = BetP_i(A_*(i)) = \sum_{\omega_j \in A_*(i)} BetP_i(\{\omega_j\}), \quad (20)$$

428 where $BetP_i$ is the predicted pignistic probability measure for pixel i . Let I_q be the
 429 set of pixels whose prediction confidence lies in the interval $(\frac{q-1}{Q}, \frac{q}{Q}]$, $q = 1, \dots, Q$.
 430 The average utility and confidence of I_q are defined, respectively, as

$$au(I_q) = \frac{1}{|I_q|} \sum_{i \in I_q} \tilde{u}_{A(i), A_*(i)}, \quad (21a)$$

431 and

$$co(I_q) = \frac{1}{|I_q|} \sum_{i \in I_q} co(i). \quad (21b)$$

432 We consider that the classifier is well calibrated if $co(I_q) \approx au(I_q)$ for all q , and
 433 we define the ECE as

$$ECE = \frac{\sum_{q=1}^Q |I_q| \times |co(I_q) - au(I_q)|}{\sum_{q=1}^Q |I_q|} \quad (22)$$

434 When only considering precise acts and labels, ECE defined by (22) boils down to
 435 the original definition in [20].

436 4.2 Precise segmentation results

437 In precise segmentation, each pixel of an image is assigned to exactly one class,
 438 the set of acts being defined as $\mathcal{F} = \{f_{\omega_1}, \dots, f_{\omega_M}\}$. Three databases without
 439 soft labels mentioned in Section 4.1 were used to train and test the E-FCNs and
 440 probabilistic FCNs (P-FCNs). The metrics defined in Section 4.1 with the utility
 441 matrix \mathbf{U} equal to the identity matrix were used for performance assessment.

442 In the experiment with each database, three widely used encoder-decoder archi-
 443 tectures were combined with the DS and utility layers, as shown in Table 4.
 444 All encoder-decoder architectures in Table 4 have the same encoder part, which
 445 consists of four stages and two convolutional layers with 3×3 kernels. Each stage
 446 is made up of three convolutional layers with 3×3 kernels and a max-pooling
 447 layer with a 2×2 non-overlapping window. Figure 4a illustrates the differences
 448 between the FCN-32s, FCN-16s, and FCN-8s architectures in their decoder parts
 449 with a deconvolutional layer. The FCN-SegNet architecture uses four deconvolu-
 450 tional layers to upsample the sparse feature maps from the end of the encoder
 451 part, as well as the feature maps from the corresponding pooling layers based on
 452 pooling indices [1], as shown in Figure 4b. The FCN-DilatedVGG architecture is
 453 the same as FCN-SegNet except that it adds a fully connected conditional ran-
 454 dom field at the end of the last deconvolutional layer [3]. The numbers P of feature
 455 maps for the Pascal, MIT and SIFT databases were, respectively, 31, 128 and 64.
 456 The numbers n of prototypes in the DS layer for these three databases were set,
 457 respectively, to 75, 300 and 95.

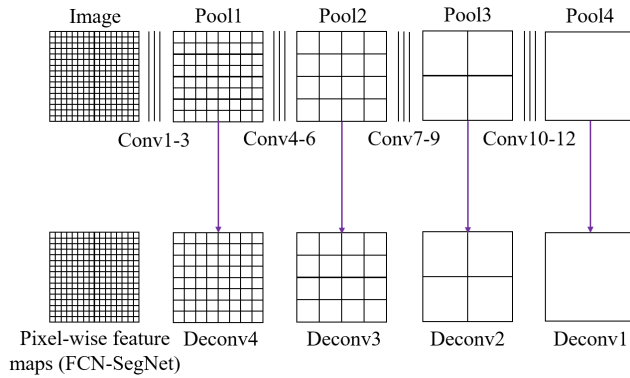
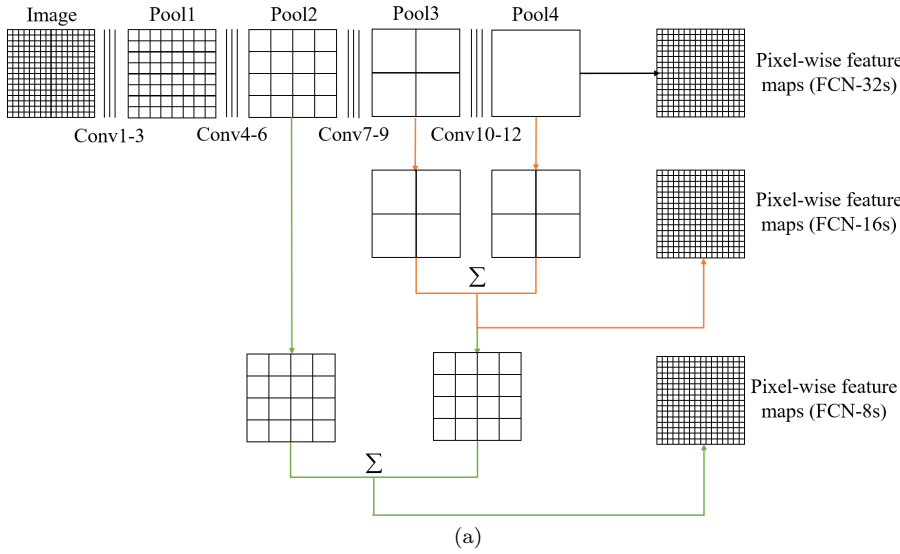


Fig. 4: Illustration of the encoder-decoder architectures used in this paper. Pooling layers are represented as grids that show relatively sparse information. Intermediate convolution layers are omitted. (a) The FCN-32s, FCN-16s, FCN-8s architectures are used to combine sparse and high-layer information with dense and low-layer information for upsampling. Black arrow: the deconvolutional layer in FCN-32s directly upsamples the outputs of Pool 4 to pixel-wise feature maps; orange arrows: the deconvolutional layer in FCN-16s combines outputs from Pool 3 and 4, lets the net predict finer details, while retaining high-level semantic information; green arrows: the deconvolutional layer in FCN-8s acquire additional feature maps from Pool 2 to provide further precision; (b) The FCN-SegNet architecture uses four deconvolutional layers to upsample the sparse feature maps from the end of the encoder part, as well as the feature maps from the corresponding pooling layers based on pooling indices (purple arrows).

Table 4: Performance evaluation of precise segmentation: (a) Pascal VOC 2011, (b) MIT-scene Parsing, and (c) SIFT Flow. P-FCN and E-FCN are, respectively, probabilistic and evidential FCNs. The rests of the notations, such as “-32s” and “-16s”, stand for different encoder-decoder architectures. **The results are in form of “mean value \pm standard deviation”.** The best results for each encoder-decoder architecture are shown in bold.

(a)		
	PU	UIoU
P-FCN-32s [25]	0.8912 \pm 0.0019	0.5941 \pm 0.0033
P-FCN-16s [25]	0.9001 \pm 0.0015	0.6243 \pm 0.0025
P-FCN-8s [25]	0.9033 \pm 0.0017	0.6269 \pm 0.0021
E-FCN-32s	0.8973 \pm 0.0021	0.6128 \pm 0.0024
E-FCN-16s	0.9045 \pm 0.0014	0.6304 \pm 0.0019
E-FCN-8s	0.9074 \pm 0.0015	0.6337 \pm 0.0020

(b)		
	PU	UIoU
P-FCN-16s [25]	0.7009 \pm 0.0030	0.289 \pm 0.0051
P-FCN-8s [25]	0.7128 \pm 0.0024	0.294 \pm 0.0048
P-FCN-SegNet [1]	0.7153 \pm 0.0023	0.305 \pm 0.0042
E-FCN-16s	0.7090 \pm 0.0026	0.292 \pm 0.0048
E-FCN-8s	0.7148 \pm 0.0025	0.296 \pm 0.0046
E-FCN-SegNet	0.7167 \pm 0.0026	0.330 \pm 0.0043

(c)		
	PU	UIoU
P-FCN-16s [25]	0.8489 \pm 0.0034	0.3922 \pm 0.0047
P-FCN-8s [25]	0.8525 \pm 0.0032	0.3948 \pm 0.0042
P-FCN-DilatedVGG [3]	0.8643 \pm 0.0036	0.4168 \pm 0.0043
E-FCN-16s	0.8521 \pm 0.0030	0.3937 \pm 0.0042
E-FCN-8s	0.8528 \pm 0.0031	0.3961 \pm 0.0040
E-FCN-DilatedVGG	0.8649 \pm 0.0035	0.4182 \pm 0.0038

458 The DS and utility layers slightly improve the accuracy of precise assign-
 459 ments performed by FCN models, **even though the performance of FCN mod-
 460 els on precise segmentation mainly depends on the encoder-decoder architecture.**
 461 Table 4a presents the results of PU and UIoU for the Pascal VOC database. E-
 462 FCNs achieved higher PU and UIoU than P-FCNs with the same encoder-decoder
 463 architecture, which shows the E-FCNs outperform the P-FCNs for precise seg-
 464 mentation. Similar improvements can also be found in the MIT-scene Parsing and
 465 SIFT Flow databases as shown, respectively, in Tables 4b and 4c.

466 The use of DS and utility layers also makes the FCN models better calibrated.
 467 Figure 5 presents a visual calibration representation of the FCN-8s models in the
 468 Pascal VOC database. The top row shows the pixel distribution of prediction
 469 confidence (21b) as histograms. The average confidence of the E-FCN-8s model
 470 closely matches its average pixel utility, while the average confidence of the P-
 471 FCN-8s model is substantially higher than its average pixel utility. This is further
 472 illustrated in the pixel utility diagrams (bottom row of Figure 5), which show
 473 pixel utility as a function of confidence. The E-FCN-8s model is well calibrated

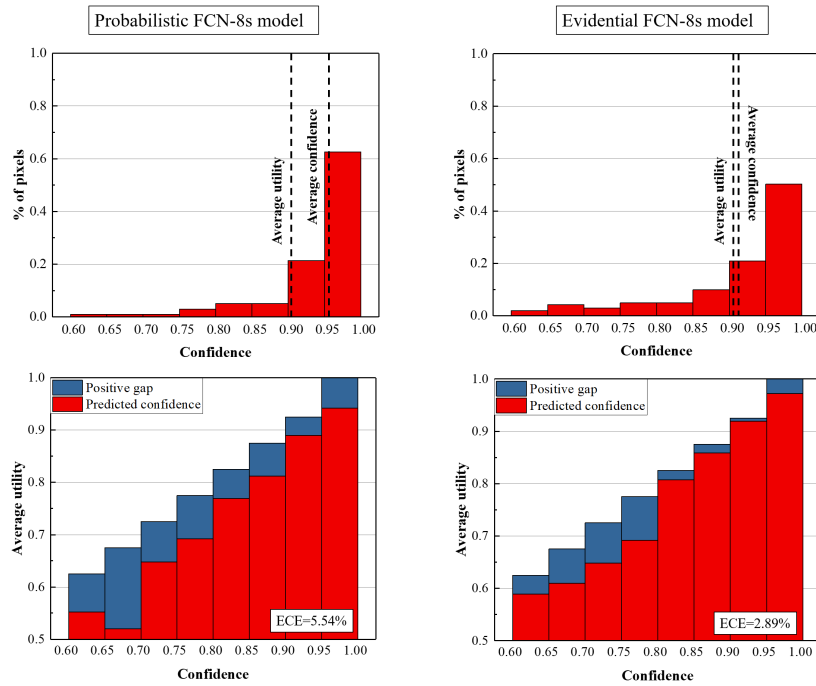


Fig. 5: Pixel confidence distributions (top) and pixel utility histograms (bottom) for P-FCN-8s (left) and E-FCN-8s (right) on the Pascal VOC database.

474 since its confidence in each bin approximates the expected average utility, whereas
 475 the predicted utility of the P-FCN-8s model does not match its confidence. As a
 476 consequence, the E-FCN-8s model achieves a smaller ECE than the probabilistic
 477 one. The effect of the DS and utility layers on the calibration can also be found
 478 in the FCN-SegNet and FCN-DilatedVGG models on the MIT-scene Parsing and
 479 SIFT Flow databases as shown, respectively, in Figures 6 and 7.

480 4.3 Imprecise segmentation results

481 In imprecise segmentation, each pixel of an image is assigned to a non-empty
 482 subset A of Ω ; the set of acts is $\mathcal{F} = \{f_A, A \in 2^\Omega \setminus \{\emptyset\}\}$, or a subset thereof. Here
 483 we only considered acts f_A such that A is a singleton, Ω or one of the soft labels
 484 in the training set. For performance evaluation, we used the metrics and the three
 485 databases described in Section 4.1. For each database, the segmentation masks
 486 with and without soft labels were used to train different FCN models. The same
 487 encoder-decoder architectures used for precise segmentation in Section 4.2 were
 488 combined with the DS and utility layers.

489 Figure 8 displays the test results according to PU and UIoU for imprecise
 490 segmentation of the Pascal VOC database. For a wide range of imprecision toler-
 491 ance degree γ , the E-FCN models reach higher PU and UIoU values than those
 492 obtained by the P-FCN models; this is due to the fact that the E-FCN models

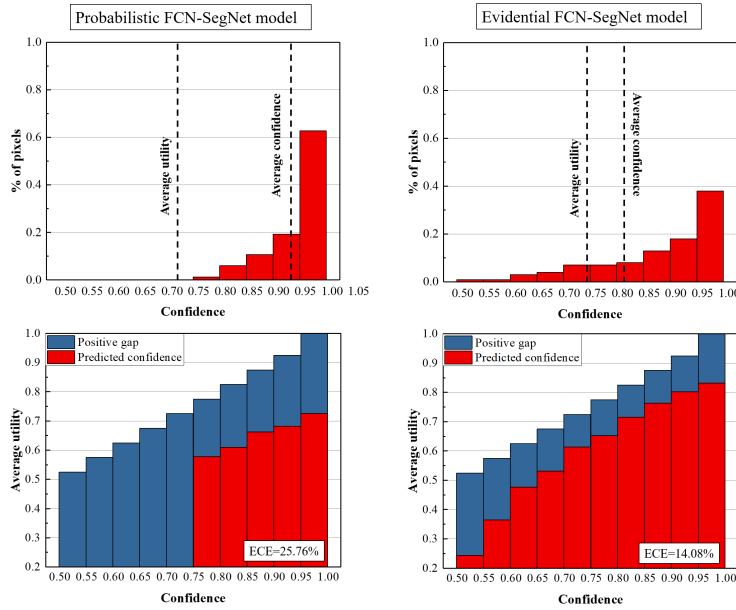


Fig. 6: Pixel confidence distributions (top) and pixel utility histograms (bottom) for P-FCN-SegNet (left) and E-FCN-SegNet (right) on the MIT-scene Parsing database.

493 tend to assign ambiguous pixels to multi-class sets, instead of making precise decisions.
 494 Such imprecise assignments avoid pixel misclassification in case of high
 495 uncertainty, especially when feature vectors from an encoder-decoder architecture
 496 do not contain sufficient information to predict a precise class, and multiple classes
 497 have similar probabilities. Figure 9 shows the pixel confidence distributions for the
 498 FCN models with $\gamma = 0.8$. We can see that the average confidences of the E-FCN
 499 models are smaller than those of the P-FCN models. This observation suggests
 500 that the E-FCN models make cautious decisions for ambiguous pixels by assigning
 501 them to multi-class sets, rather than classifying them arbitrarily into a single class.
 502 The E-FCN models are thus better calibrated than those based on P-FCN, which
 503 can be over-confident. Similar results are observed with the MIT-scene Parsing
 504 (Figures 10-11) and SIFT Flow (Figures 12-13) databases. We can thus conclude
 505 the DS and utility layers improve the performance of the FCN models in imprecise
 506 segmentation tasks by allowing the assignment of some ambiguous pixels to
 507 multi-class sets.

508 In Figures 8, 10 and 12, we can see that the value of UIoU first increases and
 509 then decreases when γ increases from 0.5 to 1. To explain this behavior, Figure
 510 14 shows some segmentation examples generated by the E-FCN-8s model trained
 511 on the Pascal VOC database with soft labels. The first and second columns of
 512 Figure 14 contain, respectively, the original images and their precise segmentation
 513 predicted masks, while the third to sixth columns show the imprecise segmentation
 514 results for values of γ ranging from 0.6 to 0.9. When γ increases from 0.5 to 0.8,
 515 the majority of the green masks (the areas whose pixels are assigned to multi-class

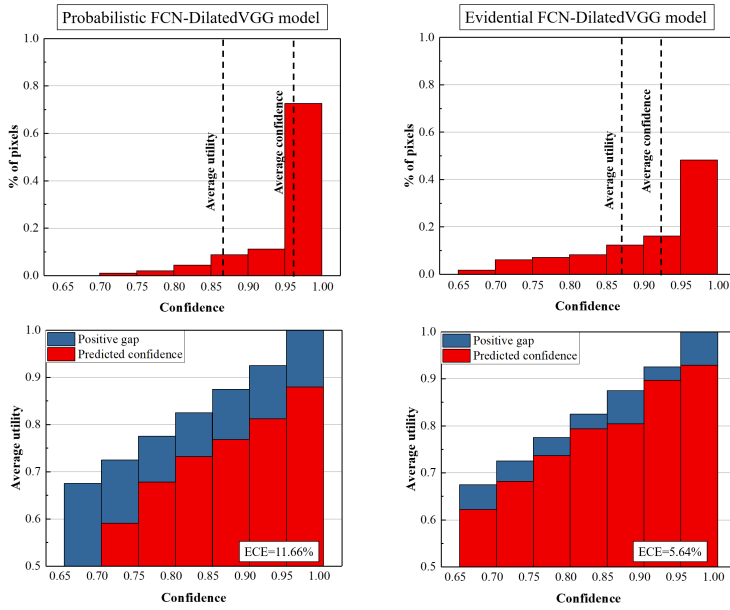


Fig. 7: Pixel confidence distributions (top) and pixel utility histograms (bottom) for P-FCN-DilatedVGG (left) and E-FCN-DilatedVGG (right) on the SIFT Flow database.

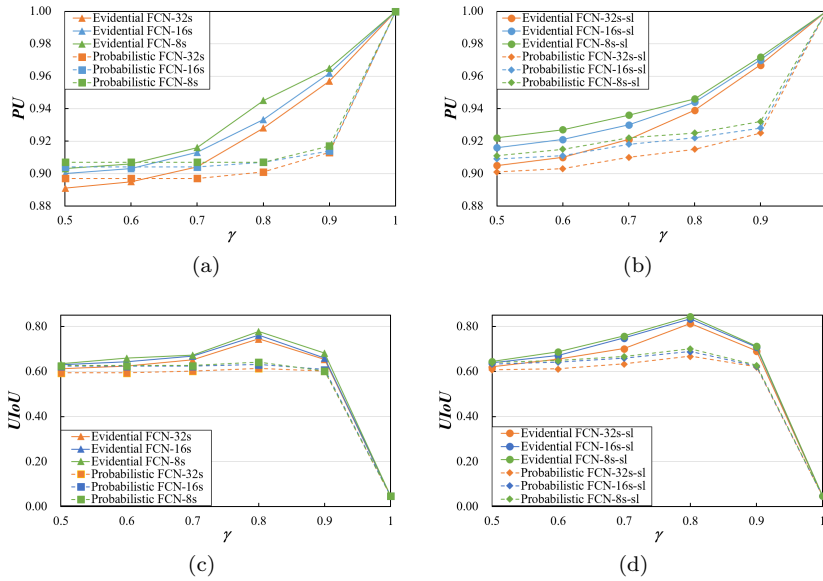


Fig. 8: Testing PU and UIoU vs. γ on the Pascal VOC database. The first and second columns are the models trained with/without soft labels, respectively.

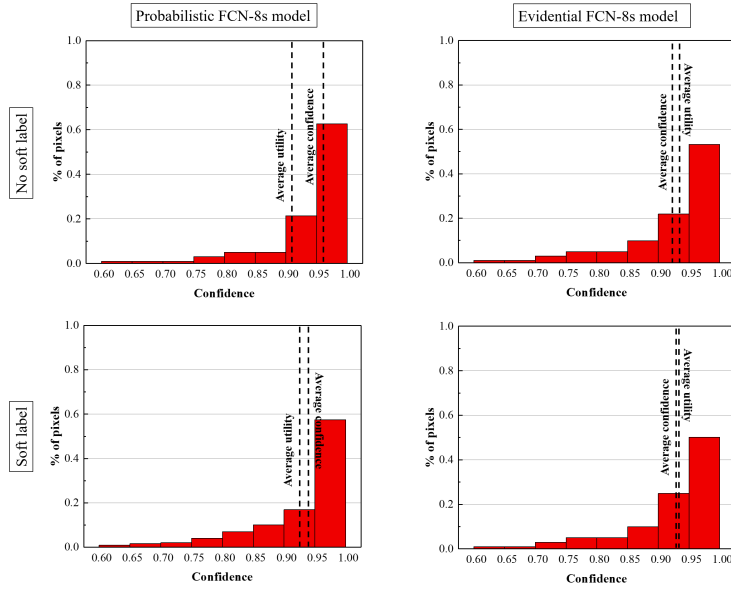


Fig. 9: Pixel confidence distributions for the P-FCN-8s (left) and E-FCN-8s (right) models on the Pascal VOC 2011 database without (top)/with (bottom) soft labels.

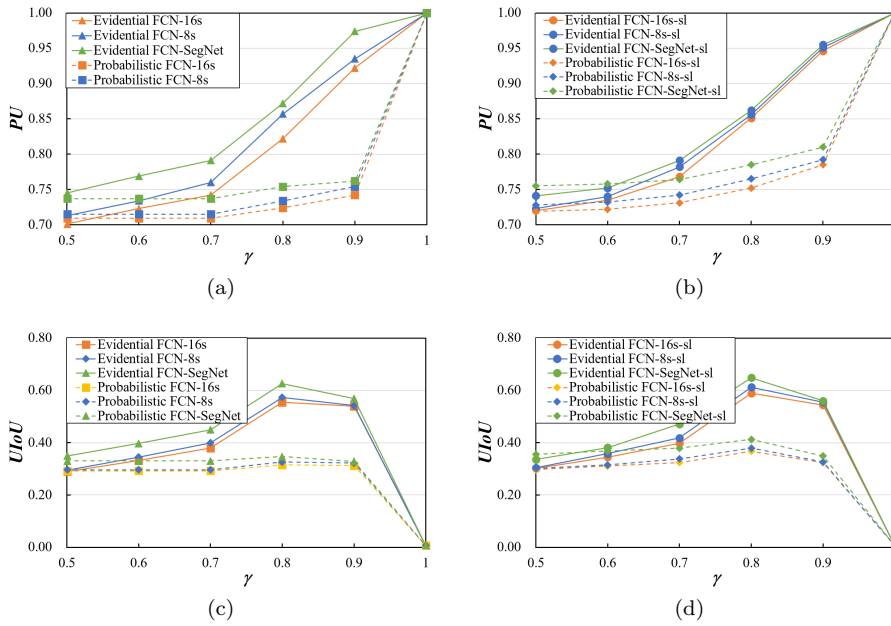


Fig. 10: Testing PU and UIoU vs. γ on the MIT-scene Parsing database. The first and second columns are the models trained with/without soft labels, respectively.

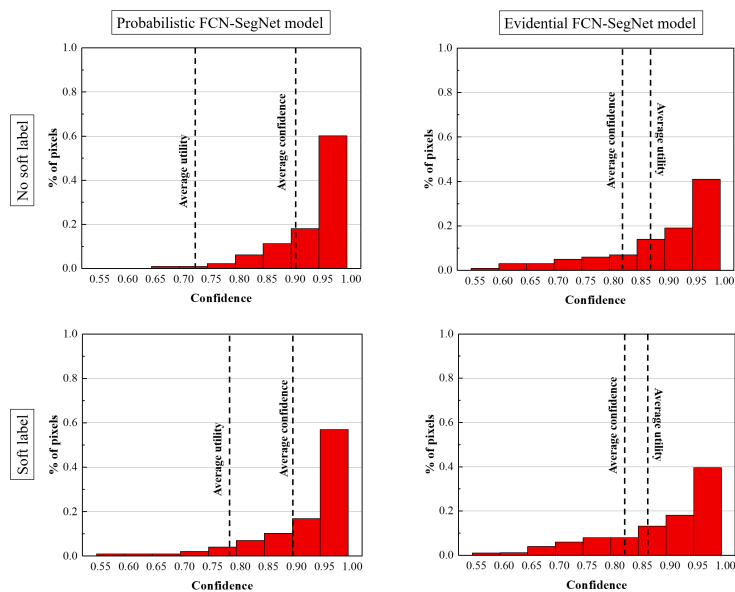


Fig. 11: Pixel rate histograms for the P-FCN-SegNet (left) and E-FCN-SegNet (right) models on the MIT-scene Parsing database without (top)/with (bottom) soft labels.

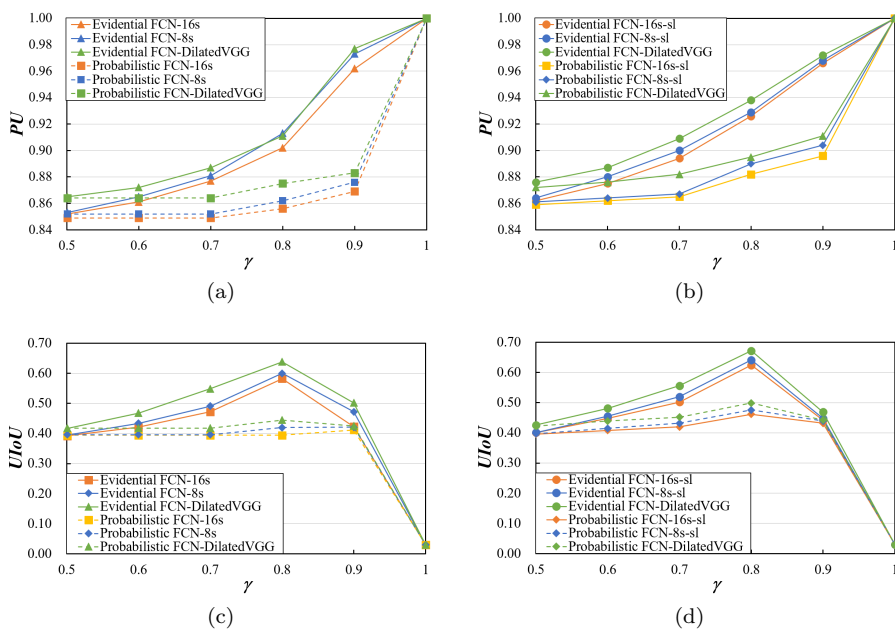


Fig. 12: Testing PU and UIoU vs. γ on the SIFT Flow database. The first and second columns are the models trained with/without soft labels, respectively.

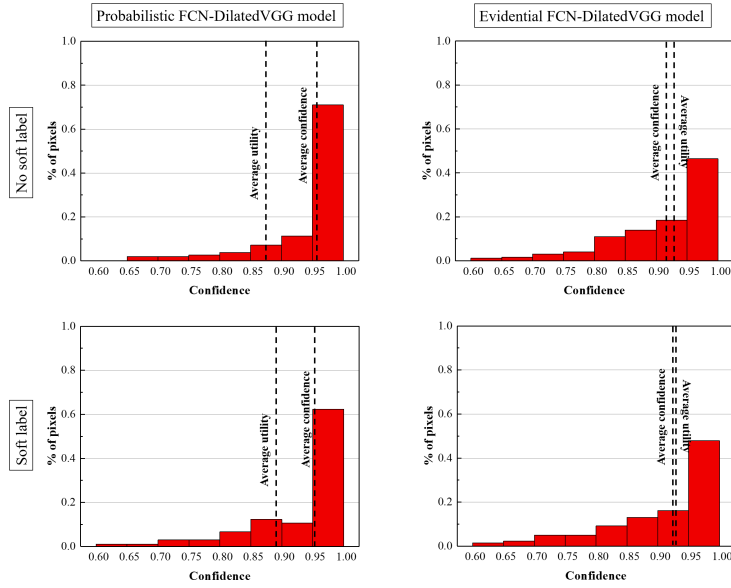


Fig. 13: Pixel rate histograms for the P-FCN-DilatedVGG (left) and E-FCN-DilatedVGG (right) models on the SIFT Flow database without (top)/with (bottom) soft labels.

sets) tends to cover the red masks (the areas whose pixels are incorrectly classified in the precise segmentation). This observation can be explained by the fact that, in Eq. (19), the increase in the utility of the intersection between predicted and labeled areas is larger than the increase in the union between the two areas. As a result, UIoU increases when γ increases from 0.5 to 0.8. However, when γ increases from 0.8 to 1.0, the majority of the green masks cover the areas predicted correctly in the precise segmentation, which causes the increase in the utility of intersection to be smaller than the increase in the union areas. This phenomenon leads to the decrease of UIoU when γ is larger than 0.8.

The use of soft labels improves the performance of the FCN models for imprecision segmentation tasks. As shown in Figure 8, the FCN models trained on the Pascal VOC database with soft labels have larger testing PU and UIoU than the ones without soft labels, which demonstrates the accuracy improvement using soft labels. Additionally, the use of soft labels can also improve the calibration of the FCN models. Figure 15 shows that the ECEs and bin gaps in the E-FCN and P-FCN models are smaller when using the learning set with soft labels. [These results demonstrate the feasibility of processing pixels with ambiguous information using soft labels when training FCN models.](#) An improvement of accuracy and calibration due to learning from soft labels is also observed with the MIT-scene Parsing and SIFT Flow databases, as shown, respectively, in Figures 16 and 17. Therefore, we can conclude that the use of soft labels improves the accuracy and calibration of FCN models.

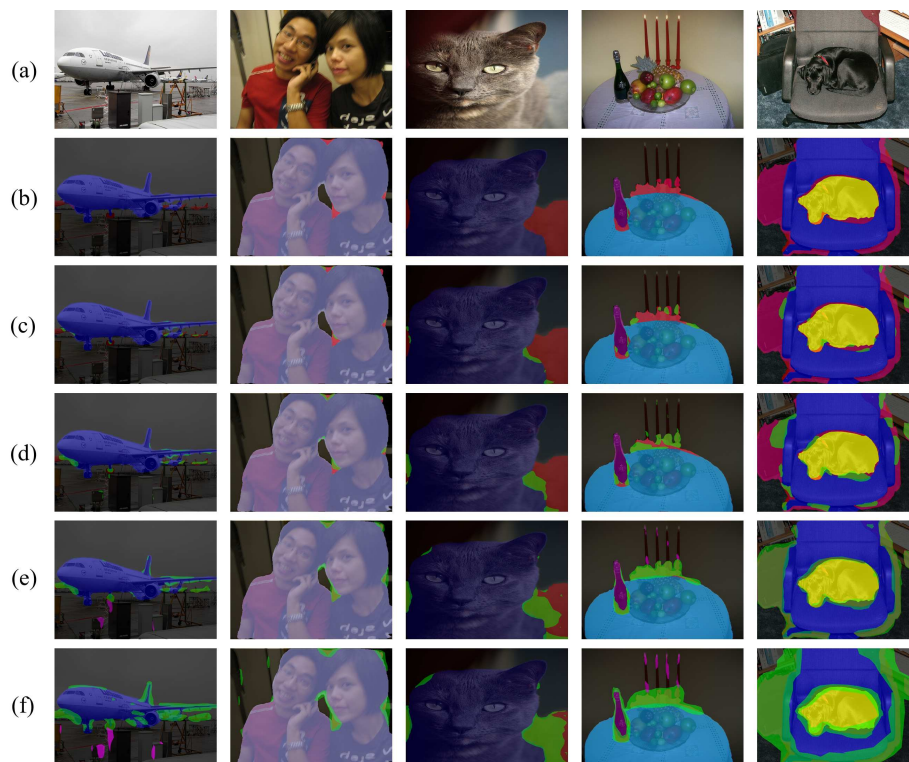


Fig. 14: Segmentation examples from the Pascal VOC 2011 database: (a) Original image, (b) Precise segmentation, (c) Imprecise segmentation with $\gamma = 0.6$, (d) Imprecise segmentation with $\gamma = 0.7$, (e) Imprecise segmentation with $\gamma = 0.8$, and (f) Imprecise segmentation with $\gamma = 0.9$. Red masks are pixels incorrectly classified in the precise segmentation; green masks are pixels assigned to multi-class sets except set Ω ; pink masks are pixels assigned to set Ω ; other masks are pixels assigned to correct single-class sets.

538 4.4 Novelty detection results

539 For novelty detection, a pixel is considered as an outlier or an ambiguous sample if
 540 it is assigned to set Ω . Figures 18, 19 and 20 show the results of novelty detection
 541 using the E-FCN and P-FCN models when the learning set is extracted, respec-
 542 tively, from the Pascal VOC, MIT-scene Parsing and SIFT Flow databases, and
 543 the test set is composed of images from the other two databases. In each testing
 544 set composed of two databases, only the pixels whose class is not represented in
 545 the corresponding learning set are reported in Figures 18-20. The E-FCN models
 546 assign outliers and some known-class pixels to set Ω for values of γ between 0.7
 547 and 0.9, while the P-FCN models do not. This observation shows that the E-FCN
 548 models are more efficient than the probabilistic ones for rejecting outliers together
 549 with ambiguous samples. [The proposed architecture thus has the potential to per-](#)
 550 [form novelty detection once given a reasonable value of tolerance to imprecision.](#)

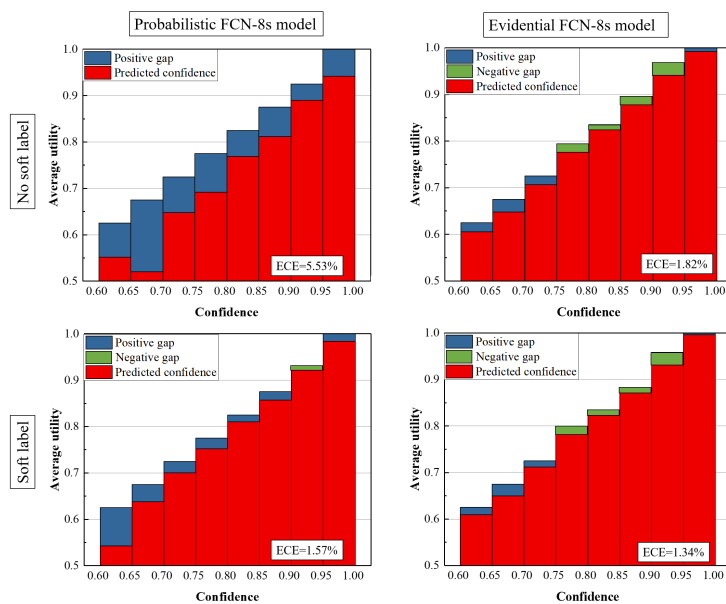


Fig. 15: Average utility histograms for P-FCN-8s (left) and E-FCN-8s (right) with $\gamma = 0.8$ on the Pascal VOC 2011 database without (top)/with (bottom) soft labels.

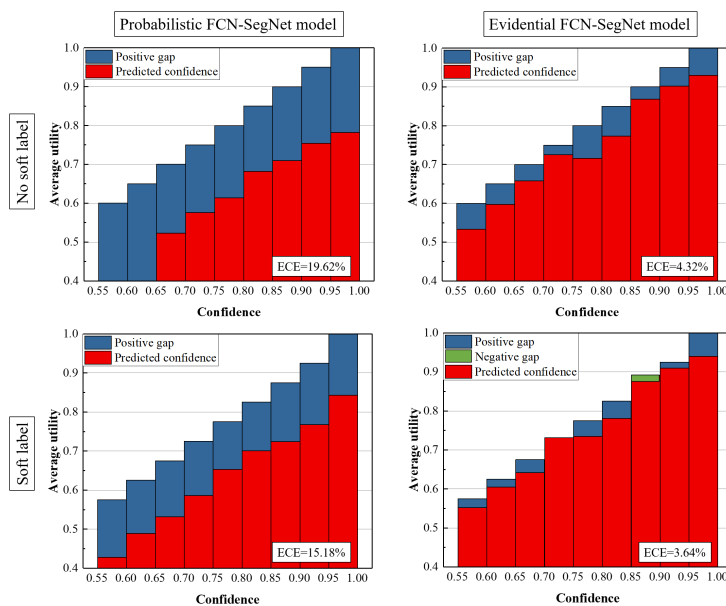


Fig. 16: Average utility histograms for P-FCN-SegNet (left) and E-FCN-SegNet (right) with $\gamma = 0.8$ on the MIT-scene Parsing database without (top)/with (bottom) soft labels.

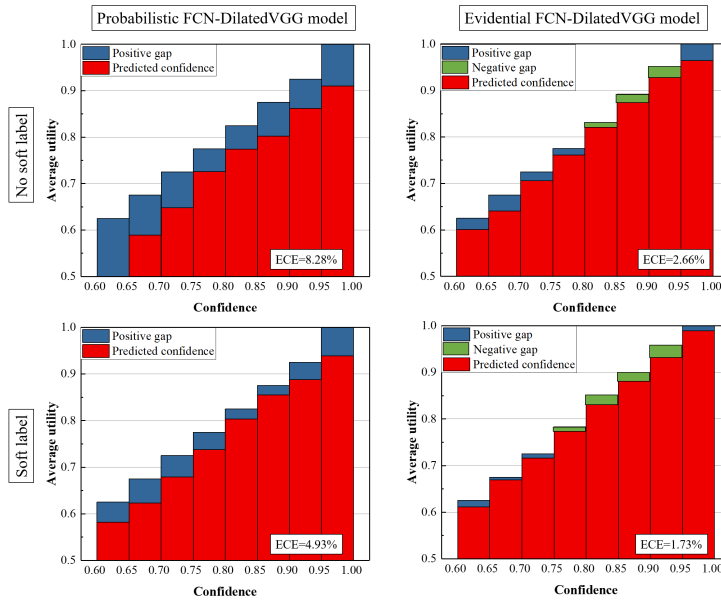


Fig. 17: Average utility histograms for P-FCN-DilatedVGG (left) and E-FCN-DilatedVGG (right) with $\gamma = 0.8$ on the SIFT Flow database without (top)/with (bottom) soft labels.

551 However, none of the FCN models performs well when γ is less than 0.7 since these
 552 models favor precise decisions.

553 The E-FCN models tend to reject unknown objects whose features are very
 554 different from those of the known objects in the learning set. For example, Figure
 555 21 shows images from the MIT-scene Parsing database in which pixels representing
 556 ‘bag’, ‘street light’ and ‘ball’ objects are rejected by an E-FCN-8s model trained
 557 using the Pascal VOC database, which does not contain these objects. As shown
 558 in Table 5, 75.2% of the pixels representing a ball in the MIT-scene Parsing and
 559 and SIFT Flow databases are assigned to Ω , while 16.1% are assigned to a set of
 560 classes containing “bottle”. For the “bag” and “street light” classes, these numbers
 561 are, respectively, 68.4%/21.8% and 77.3%/16.3%. Some unknown objects are not
 562 so easily rejected because of their similarity with known objects. For instance,
 563 84.7% of the pixels representing a seat and 81.7% of pixels representing a bench
 564 are assigned to a set of classes containing “chair”, and 88% of “wall” pixels are
 565 assigned to a set of classes containing “background”.

566 We can also observe that the FCN models trained using a learning set with
 567 soft labels reject more outliers than those trained without soft labels, as shown in
 568 Figures 18, 19 and 20. This is because the use of soft labels makes the FCN models
 569 more cautious and better calibrated, as discussed in Section 4.3. More precisely,
 570 for ambiguous pixels or outliers, the output mass functions of the FCN models
 571 trained with soft labels are more uniform than those computed by FCN models
 572 trained without soft labels. As a result, ambiguous pixels and outliers are more

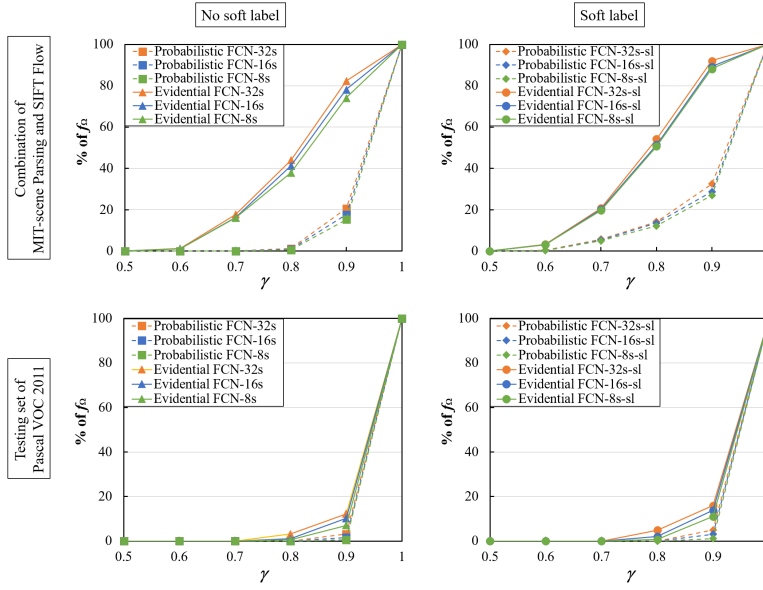


Fig. 18: Proportion of pixels assigned to Ω as a function of γ for novelty detection on the combination of MIT-scene Parsing and SIFT Flow databases (top) and the testing set from the Pascal VOC 2011 database (bottom) when the learning set is from the Pascal VOC database without (left)/with (right) soft labels.

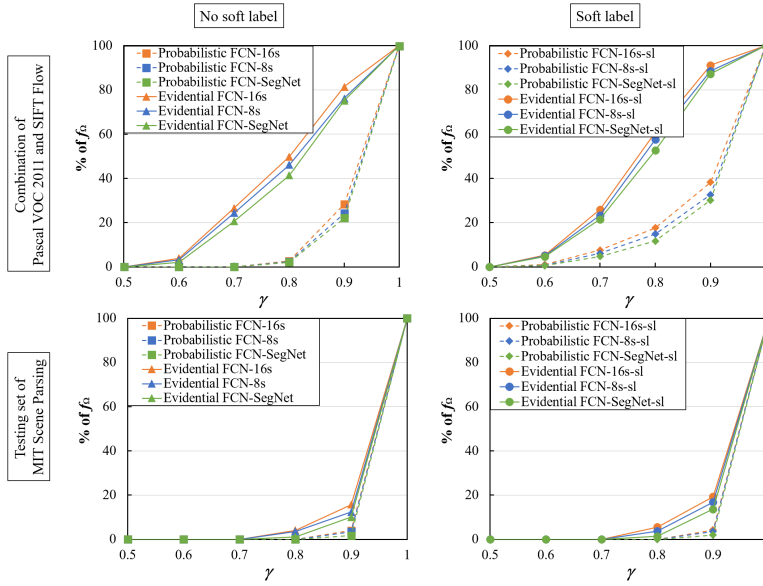


Fig. 19: Proportion of pixels assigned to Ω as a function of γ for novelty detection on the combination of Pascal VOC 2011 and SIFT Flow (top) and the testing set of the MIT-scene Parsing database (bottom) when the learning set is from the Pascal VOC database without (left)/with (right) soft labels.

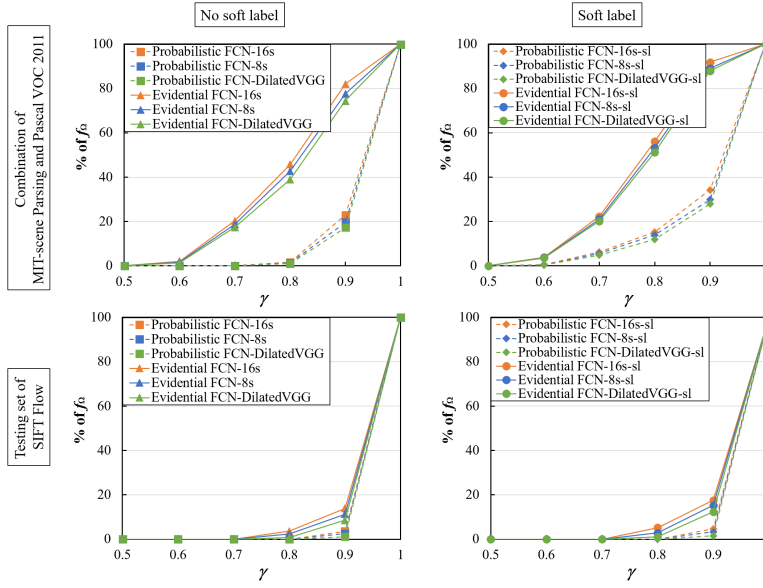


Fig. 20: Proportion of pixels assigned to Ω as a function of γ for novelty detection on the combination of Pascal VOC 2011 and MIT-scene Parsing (top) and the testing set of the SIFT Flow database (bottom) when the learning set is from the Pascal VOC database without (left)/with (right) soft labels.

Table 5: Percentage of pixels from some unknown classes in the MIT-scene Parsing and SIFT Flow databases classified by an E-FCN-8s model trained on the Pascal VOC database into some sets of classes. The model was trained with soft labels and $\gamma = 0.8$. For instance, 68.4% of the pixels representing a bag were rejected (i.e., assigned to Ω), and 84.7% of pixels representing a seat were assigned to a set of classes containing the class “chair”.

		True class						
		bag	street light	ball	seat	bench	bed	wall
Assigned set	Ω	68.4	77.3	75.2	7.8	4.7	15.9	4.9
	{bottle, ...}	21.8	16.3	16.1	48.5	39.7	30.3	0.2
	{chair, ...}	11.3	9.2	8.5	84.7	81.7	58.6	0.3
	{background, ...}	15.2	13.7	11.5	58.7	48.6	46.9	88.0
	Others	4.2	2.4	1.5	2.7	3.5	5.2	3.7

573 easily assigned to set Ω . We can thus conclude that soft labels have the potential
574 to enhance novelty detection performance.

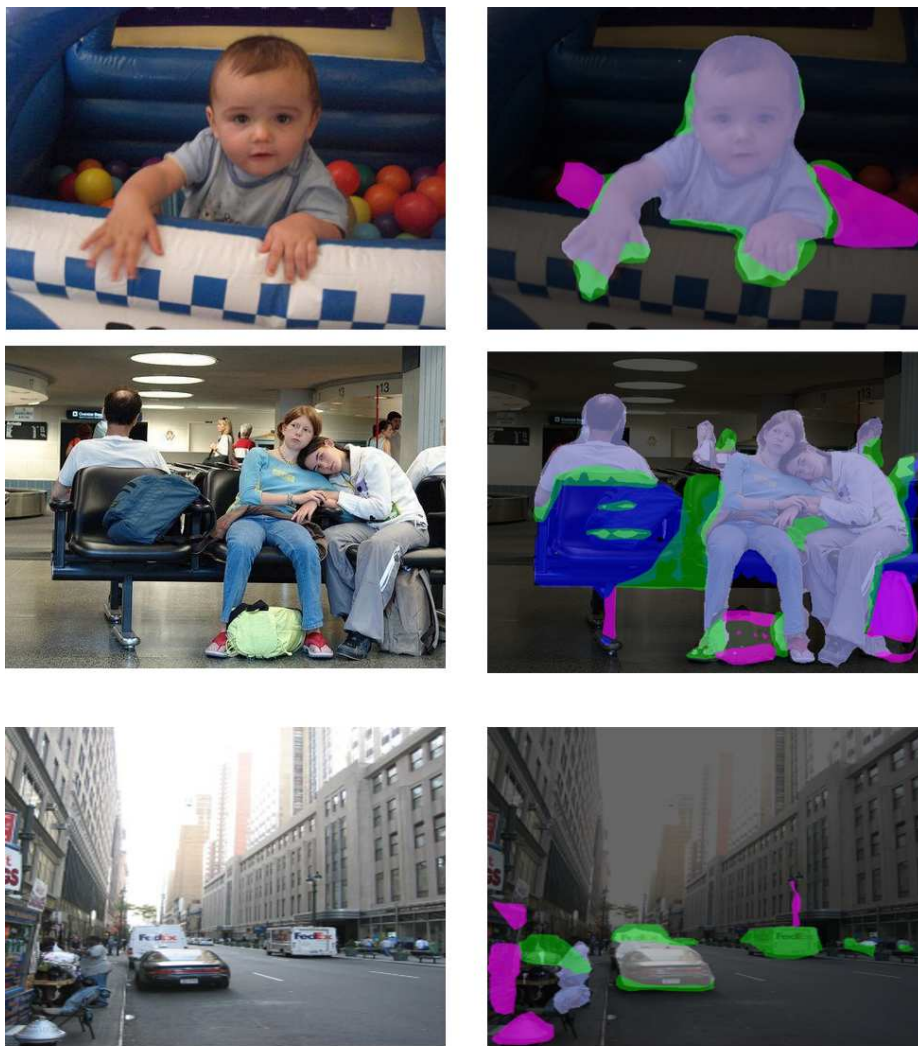


Fig. 21: Examples of novelty detection from the MIT-scene Parsing database and their segmentation masks given by the E-FCN-8s model trained using the Pascal VOC database with soft labels when γ equals 0.8. Red masks are pixels incorrectly assigned in the precise segmentation; green masks are pixels assigned to multi-class sets except set Ω ; pink masks are pixels assigned to set Ω ; other masks are pixels assigned to correct single-class sets.

575 5 Conclusions

576 In this paper, we have presented a new approach based on the combination of DS
 577 theory and FCN for image semantic segmentation. In the proposed model, called
 578 evidential fully convolutional network (E-FCN), an encoder-decoder architecture
 579 first extracts pixel-wise feature maps from an input image. A Dempster-Shafer

580 layer then computes mass functions at each pixel location based on distances to
581 prototypes. Finally, a utility layer performs semantic segmentation based on pixel-
582 wise mass functions. The proposed model can be trained using a learning set with
583 soft labels in an end-to-end way.

584 The main finding of this study is that the proposed combination of FCNs and
585 ENNs makes it possible to improve accuracy and calibration of FCN models by
586 assigning ambiguous pixels to multi-class sets, **while maintaining the good perfor-**
587 **mance of FCNs in precise segmentation tasks.** The E-FCN model is able to select
588 a set of classes when the object representation does not allow us to select a single
589 class unambiguously, which easily leads to incorrect decision-making in probabilis-
590 tic FCNs. This result provides a new direction to improve the performance of FCN
591 models for semantic segmentation. The learning strategy using soft labels further
592 improves the accuracy and calibration of the FCN models. Additionally, the pro-
593 posed approach makes it possible to reject outliers together with ambiguous pixels
594 when the tolerance to imprecision is between 0.7 and 0.9.

595 Future work will focus on two main aspects. First, we will investigate multi-
596 FCN-model information fusion for semantic segmentation based on the defini-
597 tion of soft labels, using an approach similar to that introduced in [36]. Other
598 advanced evidential classifiers, such as the contextual-discounting evidential K -
599 nearest neighbor [14] will also be considered to improve the performance of the
600 proposed neural network architecture.

601 **Acknowledgements** This research was supported by a scholarship from the China Scholar-
602 ship Council and by the Labex MS2T (reference ANR-11-IDEX-0004-02).

603 Conflict of interest

604 The authors declare that they have no conflict of interest.

605 References

- 606 1. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional
607 encoder-decoder architecture for image segmentation. *IEEE Transactions on*
608 *Pattern Analysis and Machine Intelligence* **39**(12), 2481–2495 (2017)
- 609 2. Biggio, B., Nelson, B., Laskov, P.: Support vector machines under adversarial
610 label noise. In: *Asian conference on machine learning*, pp. 97–112 (2011)
- 611 3. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab:
612 Semantic image segmentation with deep convolutional nets, atrous convolu-
613 tion, and fully connected crfs. *IEEE Transactions on Pattern Analysis and*
614 *Machine Intelligence* **40**(4), 834–848 (2017)
- 615 4. Chen, X.l., Wang, P.h., Hao, Y.s., Zhao, M.: Evidential KNN-based condi-
616 tion monitoring and early warning method with applications in power plant.
617 *Neurocomputing* **315**, 18–32 (2018)
- 618 5. Côme, E., Oukhellou, L., Dencœux, T., Aknin, P.: Learning from partially su-
619 pervised data using mixture models and belief functions. *Pattern Recognition*
620 **42**(3), 334–348 (2009)

- 621 6. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R.,
622 Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban
623 scene understanding. In: Proceedings of the IEEE conference on computer
624 vision and pattern recognition, pp. 3213–3223 (2016)
- 625 7. Dempster, A.P.: Upper and lower probabilities induced by a multivalued map-
626 ping. *Annals of Mathematical Statistics* **38**, 325–339 (1967)
- 627 8. Denœux, T.: Analysis of evidence-theoretic decision rules for pattern classifi-
628 cation. *Pattern Recognition* **30**(7), 1095–1107 (1997)
- 629 9. Denœux, T.: A neural network classifier based on Dempster-Shafer theory.
630 *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and
631 Humans* **30**(2), 131–150 (2000)
- 632 10. Denœux, T.: 40 years of Dempster-Shafer theory. *International Journal of
633 Approximate Reasoning* **79**, 1–6 (2016)
- 634 11. Denœux, T.: Decision-making with belief functions: a review. *International
635 Journal of Approximate Reasoning* **109**, 87–110 (2019)
- 636 12. Denœux, T.: Logistic regression, neural networks and Dempster-Shafer theory:
637 A new perspective. *Knowledge-Based Systems* **176**, 54–67 (2019)
- 638 13. Denœux, T., Dubois, D., Prade, H.: Representations of uncertainty in artificial
639 intelligence: Beyond probability and possibility. In: *A Guided Tour of Artificial
640 Intelligence Research*, vol. 1, chap. 4, pp. 119–150. Springer Verlag (2020)
- 641 14. Denœux, T., Kanjanatarakul, O., Sriboonchitta, S.: A new evidential k-nearest
642 neighbor rule based on contextual discounting with partially supervised learn-
643 ing. *International Journal of Approximate Reasoning* **113**, 287–302 (2019)
- 644 15. Denœux, T., Smets, P.: Classification using belief functions: relationship be-
645 tween case-based and model-based approaches. *IEEE Transactions on Sys-
646 tems, Man, and Cybernetics, Part B (Cybernetics)* **36**(6), 1395–1406 (2006)
- 647 16. Ess, A., Müller, T., Grabner, H., Van Gool, L.: Segmentation-based urban
648 traffic scene understanding. In: *BMVC*, vol. 1, p. 2. Citeseer (2009)
- 649 17. Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zis-
650 serman, A.: The pascal visual object classes challenge: A retrospective. *Inter-
651 national journal of computer vision* **111**(1), 98–136 (2015)
- 652 18. Forouzanfar, M., Forghani, N., Teshnehlab, M.: Parameter optimization of
653 improved fuzzy c-means clustering algorithm for brain mr image segmentation.
654 *Engineering Applications of Artificial Intelligence* **23**(2), 160–168 (2010)
- 655 19. Guettari, N., Capelle-Laizé, A.S., Carré, P.: Blind image steganalysis based
656 on evidential K-Nearest Neighbors. In: *Proceedings of the 2016 IEEE Interna-
657 tional Conference on Image Processing*, pp. 2742–2746. Phoenix, USA (2016)
- 658 20. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern
659 neural networks. *arXiv preprint arXiv:1706.04599* (2017)
- 660 21. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recogni-
661 tion. In: *Proceedings of the IEEE conference on computer vision and pattern
662 recognition*, pp. 770–778 (2016)
- 663 22. Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with
664 gaussian edge potentials. In: *Advances in neural information processing sys-
665 tems*, pp. 109–117 (2011)
- 666 23. Kumar, N., Berg, A.C., Belhumeur, P.N., Nayar, S.K.: Attribute and sim-
667 ilar classifiers for face verification. In: *Proceedings of the 12th International
668 Conference on Computer Vision*, pp. 365–372. IEEE, Kyoto, Japan (2009)

- 669 24. Lian, C., Ruan, S., Dencœux, T.: An evidential classifier based on feature selection and two-step classification strategy. *Pattern Recognition* **48**, 2318–2327 (2015)
- 670
- 671
- 672 25. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 3431–3440 (2015)
- 673
- 674
- 675 26. Ma, L., Dencœux, T.: Partial classification in the belief function framework. *Knowledge-Based Systems* **214**, 106742 (2021)
- 676
- 677 27. Natarajan, N., Dhillon, I.S., Ravikumar, P.K., Tewari, A.: Learning with noisy labels. In: *Advances in neural information processing systems*, pp. 1196–1204 (2013)
- 678
- 679
- 680 28. Noh, H., Hong, S., Han, B.: Learning deconvolution network for semantic segmentation. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1520–1528 (2015)
- 681
- 682
- 683 29. O’Hagan, M.: Aggregating template or rule antecedents in real-time expert systems with fuzzy set logic. In: *Twenty-Second Asilomar Conference on Signals, Systems and Computers*, vol. 2, pp. 681–689 (1988)
- 684
- 685
- 686 30. Shafer, G.: *A mathematical theory of evidence*. Princeton University Press, Princeton (1976)
- 687
- 688 31. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
- 689
- 690 32. Smets, P.: Belief functions: the disjunctive rule of combination and the generalized Bayesian theorem. *International Journal of approximate reasoning* **9**(1), 1–35 (1993)
- 691
- 692
- 693 33. Su, Z.G., Dencœux, T., Hao, Y.S., Zhao, M.: Evidential K-NN classification with enhanced performance via optimizing a class of parametric conjunctive t-rules. *Knowledge-Based Systems* **142**, 7–16 (2018)
- 694
- 695
- 696 34. Tighe, J., Lazebnik, S.: Superparsing: scalable nonparametric image parsing with superpixels. In: *European conference on computer vision*, pp. 352–365. Springer (2010)
- 697
- 698
- 699 35. Tong, Z., Xu, P., Dencœux, T.: ConvNet and Dempster-Shafer theory for object recognition. In: *Processing of the 13th international conference on Scalable Uncertainty Management*, pp. 368–381. Springer International Publishing, Cham (2019)
- 700
- 701
- 702
- 703 36. Xu, P., Davoine, F., Bordes, J.B., Zhao, H., Dencœux, T.: Multimodal information fusion for urban scene understanding. *Machine Vision and Applications* **27**(3), 331–349 (2016)
- 704
- 705
- 706 37. Yager, R.R.: On ordered weighted averaging aggregation operators in multicriteria decision-making. *IEEE Transactions on systems, Man, and Cybernetics* **18**(1), 183–190 (1988)
- 707
- 708
- 709 38. Yager, R.R., Liu, L.: *Classic works of the Dempster-Shafer theory of belief functions*, vol. 219. Springer, Berlin, Heidelberg (2008)
- 710
- 711 39. Yoon, Y., Jeon, H.G., Yoo, D., Lee, J.Y., So Kweon, I.: Learning a deep convolutional network for light-field image super-resolution. In: *Proceedings of the IEEE international conference on computer vision workshops*, pp. 24–32 (2015)
- 712
- 713
- 714
- 715 40. Yuan, B., Yue, X., Lv, Y., Dencœux, T.: Evidential deep neural networks for uncertain data classification. In: *International Conference on Knowledge Science, Engineering and Management*, pp. 427–437. Springer (2020)
- 716
- 717

-
- 718 41. Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional net-
719 works. In: European conference on computer vision, pp. 818–833. Springer
720 (2014)
- 721 42. Zeiler, M.D., Taylor, G.W., Fergus, R.: Adaptive deconvolutional networks
722 for mid and high level feature learning. In: 2011 International Conference on
723 Computer Vision, pp. 2018–2025. IEEE (2011)
- 724 43. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Se-
725 mantic understanding of scenes through the ade20k dataset. arXiv preprint
726 arXiv:1608.05442 (2016)