



UNIVERSITÀ POLITECNICA DELLE  
MARCHE

DIGITAL ADAPTIVE CIRCUITS AND LEARNING  
SYSTEMS

# Sound Event Detection con la tecnica del “*few-shot learning*”

*Matteo Orlandini e Jacopo Pagliuca*

Prof. Stefano SQUARTINI  
Dott.ssa Michela CANTARINI

15 luglio 2021

# Indice

<b>1</b>	<b>Introduzione</b>	<b>1</b>
<b>2</b>	<b>Convolutional Neural Network (CNN)</b>	<b>2</b>
<b>3</b>	<b>Area under precision-recall curve (AUPRC)</b>	<b>2</b>
<b>4</b>	<b>Few-shot learning</b>	<b>3</b>
4.1	Introduzione . . . . .	3
4.2	The meta learning framework . . . . .	4
4.3	Approcci al meta-apprendimento . . . . .	4
<b>5</b>	<b>Prototypical Network</b>	<b>5</b>
<b>6</b>	<b>Relation Network</b>	<b>7</b>
<b>7</b>	<b>Dataset Spoken Wikipedia Corpora</b>	<b>9</b>
7.1	Annotazioni . . . . .	9
7.2	Lettori . . . . .	10
7.3	Parole . . . . .	11

# 1 Introduzione

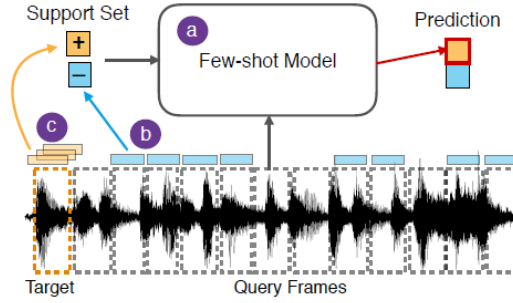


Figura 1: Metodo proposto per il few-shot sound event detection. (a) Applicazione del modello few-shot, (b) costruzione del set di esempi negativi, in blu, e (c) data augmentation per la generazione di più esempi positivi, in arancione. [4]

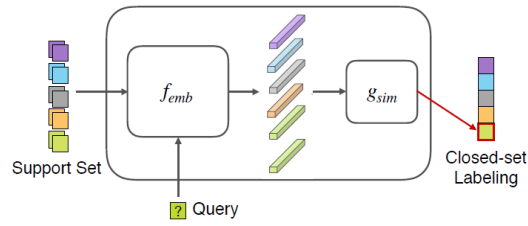


Figura 2: Modello few-shot learning nel caso 5-way 2-shot. [4]

## 2 Convolutional Neural Network (CNN)

Una rete neurale convoluzionale è una rete feedforward pensata per applicazioni che richiedono l'elaborazione di immagini o di grandi moli di dati. L'ingresso della rete è costituito da una o più matrici che attraversano delle fasi di convoluzione che ne riducono le dimensioni. I filtri di convoluzione sono detti kernel o di pooling. Infine, il risultato viene passato a una rete fully connected che ha il compito di classificazione.

La prima parte del filtro consiste nella convoluzione dell'immagine in ingresso con una serie di kernel i cui parametri vengono allenati. Il filtro viene traslato in entrambe le dimensioni della matrice producendo una serie di feature bidimensionali. Il parametro *stride* definisce il passo con cui il kernel viene traslato sulla matrice di ingresso. Solitamente, viene effettuato uno zero-padding sui contorni del volume in modo da poter caratterizzare anche i valori che si trovano ai bordi delle matrici. In seguito alla convoluzione viene applicata una funzione non-lineare come ad esempio una sigmoide o una ReLu con lo scopo di aumentare la proprietà di non linearità.

Il pooling layer è uno strato della rete che ha lo scopo di ridurre le dimensioni della matrice prodotta dal convolutional layer. Combinando gruppi di elementi della matrice (solitamente 2x2) restituisce il valore massimo tra essi, che andrà a sostituire il blocco stesso.

## 3 Area under precision-recall curve (AUPRC)

Nella fase di test, la rete elabora una predizione dell'output a partire da un ingresso noto che poi viene confrontata con il valore effettivo. Nel nostro caso è presente un positive set e un negative set, si tratta quindi di classificazione binaria. Il confronto tra predizione e label può produrre quattro risultati:

- Vero Negativo (TN): il valore reale è negativo e il valore predetto è negativo;
- Vero Positivo (TP): il valore reale è positivo e il valore predetto è positivo;
- Falso Negativo (FN): il valore reale è positivo e il valore predetto è negativo;
- Falso Positivo (FP): il valore reale è negativo e il valore predetto è positivo;

Questi valori vanno a comporre la confusion matrix. Nel nostro caso siamo interessati a Precision e Recall.

Il parametro Precision rappresenta quanti tra i casi predetti come positivi sono realmente positivi.

Il parametro Recall indica quanti tra i casi realmente positivi è stato predetto in modo corretto.

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

## 4 Few-shot learning

(<https://www.borealisai.com/en/blog/tutorial-2-few-shot-learning-and-meta-learning-i/>)

### 4.1 Introduzione

Gli esseri umani possono riconoscere nuove classi di oggetti partendo da pochissimi esempi. Tuttavia, la maggior parte delle tecniche di machine learning richiedono migliaia di esempi per ottenere prestazioni simili a quelle umane. L'obiettivo del *few-shot learning* è classificare i nuovi dati dopo aver visto solo pochi esempi di training. Nel caso estremo, potrebbe esserci solo un singolo esempio per ogni classe (*one shot learning*). In pratica, il few-shot learning è utile quando è difficile trovare esempi di training (ad es. casi di una malattia rara) o quando il costo dell'etichettatura dei dati è elevato.

L'apprendimento few-shot viene solitamente studiato utilizzando la classificazione *N-way-K-shot*. L'obiettivo è quello di discriminare le N classi composte da K esempi ciascuna. Una tipica dimensione del problema potrebbe essere quella di discriminare tra  $N = 10$  classi con solo  $K = 5$  campioni ciascuno di training. Non possiamo allenare un classificatore usando metodi convenzionali; qualsiasi algoritmo di classificazione moderno dipenderà da molti più parametri rispetto agli esempi di addestramento e generalizzerà male.

Se i dati non sono sufficienti per ridurre il problema, una possibile soluzione è acquisire esperienza da altri problemi simili. A tal fine, la maggior parte degli approcci caratterizza l'apprendimento a breve termine con un problema di meta-apprendimento.

## 4.2 The meta learning framework

Nel framework di apprendimento classico, impariamo come classificare dai dati di training e valutiamo i risultati utilizzando i dati di test. Nel quadro del meta-apprendimento, *impariamo come imparare* a classificare in base a una serie di *training task* e valutiamo utilizzando una serie di test task; In altre parole, usiamo un insieme di problemi di classificazione per aiutare a risolvere altri insiemi non correlati.

Qui, ogni attività imita lo scenario few-shot, quindi per la classificazione N-way-K-shot, ogni attività include N classi con K esempi ciascuna. Questi sono noti come support set per il task e vengono utilizzati per apprendere come risolvere il task. Inoltre, esistono ulteriori esempi delle stesse classi, note come set di query, utilizzate per valutare le prestazioni del task corrente. Ogni task può essere completamente unico; potremmo non vedere mai le classi di un task in nessuno degli altri. L'idea è che il sistema veda ripetutamente istanze (task) durante l'addestramento che corrispondono alla struttura dell'attività finale di few-shot, ma che contengono classi diverse.

Ad ogni fase del meta-apprendimento, aggiorniamo i parametri del modello in base ad un training task selezionato casualmente. La funzione di loss è determinata dalle prestazioni di classificazione sul set di query del task, in base alla conoscenza acquisita dal relativo set di supporto. Poiché la rete viene sottoposta ad un compito diverso in ogni fase temporale, deve imparare a discriminare le classi di dati in generale, piuttosto che un particolare sottoinsieme di classi.

Per valutare le prestazioni in pochi colpi, utilizziamo una serie di attività di test. Ciascuno contiene solo classi invisibili che non erano in nessuna delle attività di formazione. Per ciascuno, misuriamo le prestazioni sul set di query in base alla conoscenza del loro set di supporto.

## 4.3 Approcci al meta-apprendimento

Gli approcci al meta-apprendimento sono diversi e non c'è consenso sull'approccio migliore. Tuttavia, esistono tre famiglie distinte, ognuna delle quali sfrutta un diverso tipo di conoscenza a priori:

- Conoscenze a priori sulla somiglianza: apprendiamo degli embedding nei training task che tendono a separare classi diverse anche quando non sono visibili.
- Conoscenze a priori sull'apprendimento: utilizziamo le conoscenze a priori per vincolare l'algoritmo di apprendimento a scegliere parametri che generalizzino bene da pochi esempi.

- Conoscenza a priori dei dati: sfruttiamo le conoscenze a priori sulla struttura e la variabilità dei dati e questo ci consente di apprendere modelli praticabili da pochi esempi.

## 5 Prototypical Network

Nella classificazione few-shot per ogni episodio viene fornito un support set di  $N$  esempi etichettati  $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  dove  $\mathbf{x}_i \in \mathbb{R}^D$  rappresenta il vettore  $D$ -dimensionale della feature (nel nostro caso spettrogrammi di dimensione  $128 \times 51$ ) e  $y_i \in \{1, \dots, K\}$  la rispettiva label.  $S_k$  denota il set di esempi etichettati con la classe  $k$ .

La rete Prototypical calcola una rappresentazione  $M$ -dimensionale  $\mathbf{c}_k$ , o *prototipo*, di ogni classe tramite una funzione di embedding  $f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$  con parametri da allenare  $\phi$ . La funzione di embedding è rappresentata da una rete convoluzionale. Ogni prototipo è il vettore media tra gli embedding delle istanze della stessa classe.

$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}_i, y_i) \in S_k} f_\phi(\mathbf{x}_i) \quad (3)$$

Data una funzione distanza  $d : \mathbb{R}^M \times \mathbb{R}^M \rightarrow [0, +\infty)$ , la rete Prototypical calcola la relazione di una query  $\mathbf{x}$  rispetto ai prototipi tramite la funzione softmax delle distanze prese con segno negativo.

$$p_\phi(y = k|\mathbf{x}) = \frac{\exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k))}{\sum_k' \exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k'))} \quad (4)$$

Il processo di training avviene minimizzando il negativo del logaritmo della probabilità  $J(\phi) = -\log(p_\phi(y = k|\mathbf{x}))$  considerando la distanza fra query e il prototipo della sua classe. Gli episodi di training sono formati campionando  $C$  classi di parole da un lettore e selezionando per ognuna casualmente  $K$  istanze e  $Q$  query.

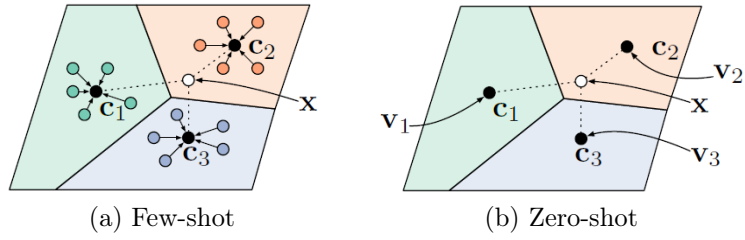


Figura 3: Reti Prototypical negli scenari few-shot e zero-shot. (a): i prototipi few-shot  $\mathbf{c}_k$  sono calcolati come la media degli embedding del support set per ogni classe. (b): i prototipi zero-shot  $\mathbf{c}_k$  sono prodotti facendo l'embedding dei metadati  $\mathbf{v}_k$ . In entrambi i casi i punti degli embedding delle query sono classificati facendo il softmax sulle distanze del prototipo delle classi:  $p_\phi(y = k|\mathbf{x}) \propto \exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k))$ . [2]



## 6 Relation Network

La Relation Network è costituita da due moduli: un modulo di *embedding*  $f_\varphi$  (equivalente a quello nella Prototypical) e un modulo di *relation*  $g_\phi$ , come illustrato in figura 4. Le istanze  $x_i$  del query set  $\mathcal{Q}$  e quelle  $x_j$  del support set  $\mathcal{S}$  vengono date in ingresso al modulo di embedding producendo dei vettori (feature maps)  $f_\varphi(x_i)$  e  $f_\varphi(x_j)$ . Questi ultimi vengono poi dati all'operatore  $\mathcal{C}(\cdot, \cdot)$  che ne fa la concatenazione:  $\mathcal{C}(f_\varphi(x_i), f_\varphi(x_j))$ . Le feature map concatenate passano poi attraverso il modulo di decisione che restituisce uno scalare da 0 a 1, il quale rappresenta la somiglianza tra  $x_i$  e  $x_j$ . Per il caso  $C$ -way one-shot, viene concatenata la query con le istanze delle  $C$  classi producendo  $C$  punteggi di somiglianza.

$$r_{i,j} = g_\phi(\mathcal{C}(f_\varphi(x_i), f_\varphi(x_j))), \quad i = 1, 2, \dots, C \quad (5)$$

Nel caso  $C$ -way K-shot invece, la query viene concatenata con la somma elemento per elemento degli embedding di ogni istanza delle classi. Quindi, in entrambi i casi i confronti  $r_{i,j}$  sono  $C$  per ogni query.

Per allenare il modello viene usato l'errore quadratico medio (MSE) in modo che l'uscita del modulo di decisione produca 1 se i vettori concatenati sono della stessa classe e 0 altrimenti.

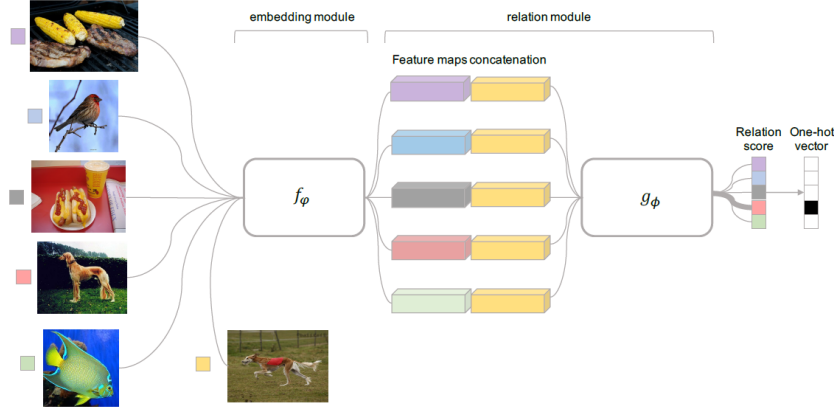


Figura 4: Architettura della Relation Network nel caso 5-way 1-shot con un esempio di query. [3]

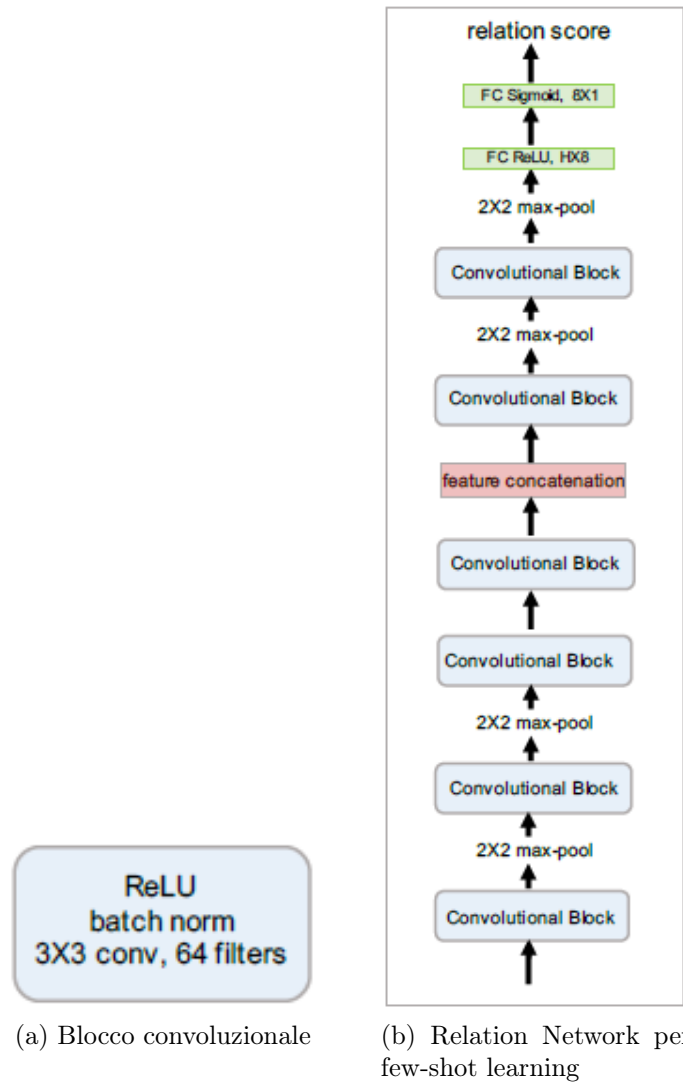


Figura 5: Architettura della Relation Network per few-shot learning (b) composta dagli elementi inclusi nel blocco convoluzionale (a). [3]

## 7 Dataset Spoken Wikipedia Corpora

Il progetto Spoken Wikipedia unisce lettori volontari di articoli di Wikipedia. Centinaia di articoli in inglese, tedesco e olandese sono disponibili per gli utenti che non sono in grado o non vogliono leggere la versione scritta dell'articolo. Il dataset trasforma i file audio in un corpus allineato nel tempo, rendendolo accessibile per la ricerca.[1]

Questo corpus ha diverse caratteristiche importanti:

- centinaia di ore di audio allineato
- lettori eterogenei
- diversi argomenti
- genere testuale ben studiato
- le annotazioni delle parole possono essere mappate all'html originale
- allineamenti a livello di fonema

Ogni articolo è suddiviso in sezioni, frasi e token. Ogni token è normalizzato e la normalizzazione è allineata all'audio.

INSERIRE SCHEMA SWC.RNC!!

### 7.1 Annotazioni

```
<article>
<meta>
<link key="DC.conformsto" value="http://nats.gitlab.io/swc/schema/swc-1.0.rnc"/>
<prop key="DC.creator" value="Spoken Wikipedia Corpus Collection Software"/>
<prop key="DC.publisher" value="Universität Hamburg"/>
<link key="DC.reference" value="http://nbn-resolving.de/urn:nbn:de:gbv:18-228-7-2209"/>
<prop key="DC.type" value="dataset"/>
<prop key="DC.license" value="CC-BY-SA"/>
<prop key="DC.title" value="Limerence"/>
<prop key="DC.language" value="en"/>
<prop key="DC.identifier" value="Limerence"/>
<prop key="DC.date.read" value="2005-04-29 00:00:00"/>
<link key="DC.source" value="https://en.wikipedia.org/wiki/Limerence"/>
<prop key="DC.source.wikiID" value="154147"/>
<prop key="DC.source.revision" value="791627969"/>
<link key="DC.source.text" value="https://en.wikipedia.org/w/index.php?title=Limerence&oldid=13811989"/>
<link key="DC.source.audio" value="https://upload.wikimedia.org/wikipedia/commons/aa/Limerence.ogg" group="audio1"/>
<prop key="DC.source.audio.offset" value="0.0" group="audio1"/>
<link key="DC.source.audio.page" value="https://en.wikipedia.org/w/index.php?title=File%3aLimerence.ogg" group="audio1"/>
<link key="DC.source.audio.date" value="2007-09-14 19:23:05" group="audio1"/>
<prop key="reader.name" value="the Eopt"/>
<prop key="processing.step" value="tokenize" group="tokenize"/>
<prop key="processing.step.date" value="2017-08-10T07:44:44.095+02:00[Europe/Berlin]" group="tokenize"/>
<prop key="processing.step.options" value="Namespace(output=articles/Limerence/tokenized.swc, all_sections=false, null_normalize=false, raw_output=null, subparser_name=tokenize, lang=en, no_introduction=false, article_dir=articles/Limerence)" group="tokenize"/>
```

```

<prop key="processing.step.git.commit.id" value="7431
dbf93f212ad828208abaf8f518fb8de11ff3" group="tokenize"/>
<prop key="processing.step.git.commit.time" value="09.08.2017 @ 15:21:55 CEST" group
="tokenize"/>
<prop key="processing.step" value="align" group="align"/>
<prop key="processing.step.date" value="2017-08-11T16:12:18.423+02:00[Europe/Berlin]
" group="align"/>
<prop key="processing.step.options" value="Namespace(output=articles/Limerence/
aligned.swc, transcript=articles/Limerence/tokenized.swc, g2p=../model_en/model.
fst.ser, phone=false, subparser_name=align, dict=../model_en/empty.dic,
acoustic_model=../model_en/, audio=articles/Limerence/audio.wav)" group="align"
/>
<prop key="processing.step.git.commit.id" value="7431
dbf93f212ad828208abaf8f518fb8de11ff3" group="align"/>
<prop key="processing.step.git.commit.time" value="09.08.2017 @ 15:21:55 CEST" group
="align"/>
</meta>

```

Codice 1: Metadati delle annotazioni delle parole in un audio

```

<d>
<extra text="LimerenceFrom wikipedia, the free encyclopedia at e n dot wikipedia dot
org.">
<s text="LimerenceFrom wikipedia, the free encyclopedia at e n dot wikipedia dot org
.">
<t text="LimerenceFrom">
<n pronunciation="LimerenceFrom" start="140" end="1190"/>
</t>
<t text="wikipedia">
<n pronunciation="wikipedia" start="1250" end="1950"/>
</t>
<t text=","/>
<t text="the">
<n pronunciation="the" start="1950" end="2070"/>
</t>
<t text="free">
<n pronunciation="free" start="2070" end="2300"/>
</t>
<t text="encyclopedia">
<n pronunciation="encyclopedia" start="2300" end="3220"/>
</t>
<t text="at">
<n pronunciation="at"/>
</t>
<t text="e">
<n pronunciation="e" start="3490" end="3710"/>
</t>

```

Codice 2: Annotazioni delle parole in un audio

## 7.2 Lettori

```

1 import json
2
3 def write_json_file(filename, list_of_dict, indent = 0):
4     f = open(filename, "w")
5     f.write(json.dumps(list_of_dict, indent = indent))
6     f.close
7
8 def read_json_file(filename):
9     f = open(filename, "r")
10    list_of_dict = json.load(f)
11    f.close
12    return list_of_dict

```

Codice 3: json\_manager.py

```

{
  "reader_name": "the epopt",
  "folder": [
    "(I_Can%27t_Get_No)_Satisfaction",
    "Ceremonial_ship_launching",
    "Limerence",
    "Revolt_of_the_Admirals",
    "Ship_commissioning"
  ]
},
{
  "reader_name": "wodup",
  "folder": [
    "0.999..%2e",
    "Execution_by_elephant",
    "Hell_Is_Other_Robots",
    "Tom_Bosley",
    "Truthiness"
  ]
},

```

Codice 4: Formato JSON dei lettori

## 7.3 Parole

```

[
  {
    "word": "i",
    "frequency": 12,
    "start": [
      660,
      8800,
      115050,
      116300,
      117910,
      228560,
      273900,
      497150,
      518270,
      534740,
      543420,
      589280
    ],
    "end": [
      870,
      8940,
      115240,
      116360,
      118020,
      228690,
      274000,
      497340,
      518520,
      534860,
      543700,
      589360
    ]
  },

```

---

Codice 5: Formato JSON delle parole allineate

## Riferimenti bibliografici

- [1] Arne Köhn, Florian Stegen e Timo Baumann. «Mining the Spoken Wikipedia for Speech Data and Beyond». Inglese. In: *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)* (23–28 mag. 2016). A cura di Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk e Stelios Piperidis. Portorož, Slovenia: European Language Resources Association (ELRA). ISBN: 978-2-9517408-9-1.
- [2] Jake Snell, Kevin Swersky e Richard S. Zemel. *Prototypical Networks for Few-shot Learning*. 2017. arXiv: 1703.05175 [cs.LG].
- [3] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip H. S. Torr e Timothy M. Hospedales. «Learning to Compare: Relation Network for Few-Shot Learning». In: *CoRR* abs/1711.06025 (2017). arXiv: 1711.06025. URL: <http://arxiv.org/abs/1711.06025>.
- [4] Yu Wang, Justin Salamon, Nicholas J. Bryan e Juan Pablo Bello. *Few-Shot Sound Event Detection*. 2020, pp. 81–85. DOI: 10.1109/ICASSP40776.2020.9054708.