

# Single-Image Crowd Counting via Multi-Column Convolutional Neural Network

## What is the problem?

The problem of single-image crowd counting involves accurately estimating the number of people in crowded scenes from images, which becomes particularly challenging in dense regions or environments where objects blend into the background, such as underwater scenes. Traditional methods, which often rely on supervised learning approaches like density maps and regression models, require extensive manual annotations, such as point-level labels for each object. This process is time-consuming, costly, and difficult, especially in complex, dense scenes. Moreover, in indiscernible object counting (IOC) tasks, such as counting fish in underwater environments, the challenge is compounded by the difficulty in distinguishing objects from their surroundings. Previous research in crowd counting has focused on supervised and semi-supervised methods, while indiscernible object counting is a relatively new area, with few existing datasets or approaches

## What has been done earlier?

Earlier work in crowd counting relied heavily on supervised learning techniques, where point-level annotations (marking each individual head) were required to generate density maps. These methods, while accurate, demanded extensive manual labeling, which is both expensive and labor-intensive. To reduce this burden, semi-supervised and weakly-supervised approaches were introduced, which require less labeled data but still depend on some degree of annotation. In the context of indiscernible object counting (IOC), existing approaches like generic object counting (GOC) and dense object counting (DOC) have made progress, but there is little focus on objects that blend into their environments.

Generic Object Counting



Person: 8 Bus: 1 Bicycle: 1

Dense Object Counting



Person: 118

Indiscernible Object Counting



Fish: 20



Fish: 61

What are the remaining challenges? What novel solution proposed by the authors to solve the problem?

The remaining challenges in crowd counting and IOC revolve around the high costs of annotation, the difficulty of detecting objects in visually ambiguous or indiscernible scenes, and the need for methods that generalize well across diverse environments. In response to these challenges, recent advances have introduced novel solutions, such as CrowdCLIP and IOCFormer, which leverage unsupervised and hybrid approaches to improve accuracy while reducing the need for manual intervention. CrowdCLIP tackles the crowd counting problem by utilizing a vision-language model in an unsupervised framework, bypassing the need for labeled data altogether. It introduces a ranking-based contrastive learning method to fine-tune the image encoder and employs a progressive filtering strategy to isolate the most relevant crowd patches, allowing for more accurate counts without annotation. Similarly, IOCFormer addresses the issue of indiscernible object counting by combining density-based and regression-based approaches in a unified framework. This method is particularly effective in counting objects that blend into their environment, as it estimates object density while directly regressing object locations. Both approaches represent significant advancements in reducing the dependency on manual labels and enhancing performance in complex, densely packed, or visually ambiguous environments.

Abhay Rajpoot, B421003