
Détection d'animaux camouflés par information temporelle et modèles Vision-Langage

Mounir Ammam

Université de Montréal

mounir.ammam@umontreal.ca

Byungsuk Min

Université de Montréal

byungsuk.min@umontreal.ca

Yanis Chikhar

Université de Montréal

yanis.chikhar@umontreal.ca

1 Revue de la littérature

2 La détection d'objets camouflés constitue un problème particulièrement difficile en vision par
3 ordinateur, car les objets d'intérêt présentent une forte similarité visuelle avec leur environnement.
4 Dans de nombreux cas, les textures, couleurs et contours de l'objet se confondent avec le fond, rendant
5 l'identification ambiguë, même pour un observateur humain. Les approches basées uniquement sur
6 des images statiques sont donc limitées lorsque les indices visuels sont faibles.

7 1.1 Limitations de la détection à partir d'images statiques

8 Li et al. (1) mettent clairement en évidence cette limitation dans le contexte des animaux camouflés.
9 Les auteurs montrent que ces derniers sont presque indiscernables du fond dans des images isolées et
10 deviennent détectables principalement lorsqu'ils sont en mouvement. Cette observation souligne que
11 l'information contenue dans une seule image est souvent insuffisante pour résoudre le problème du
12 camouflage, et motive l'exploration de sources d'information supplémentaires.

13 1.2 Exploitation du mouvement et de l'information temporelle

14 Afin de dépasser les limites des approches basées uniquement sur l'image, plusieurs travaux ont
15 étudié l'exploitation du mouvement et de l'information temporelle dans les vidéos. Le mouvement
16 constitue un indice discriminant important, car il permet de révéler des objets camouflés qui restent
17 invisibles dans des images statiques.

18 Cheng et al. (2) s'intéressent à la détection d'objets camouflés dans des séquences vidéo et montrent
19 que la prise en compte de l'information temporelle améliore significativement les performances de
20 détection. Leur travail met en évidence que les variations temporelles, même subtiles, fournissent des
21 signaux utiles pour distinguer les objets camouflés de leur arrière-plan. Ces résultats confirmont que
22 l'exploitation du mouvement est une direction prometteuse pour la détection d'animaux camouflés.

23 1.3 Vision-Language Models et connaissances sémantiques

24 En parallèle, les modèles Vision–Langage (VLM) ont récemment attiré beaucoup d'attention en
25 raison de leur capacité à apprendre des représentations conjointes image–texte à grande échelle. Des
26 modèles tels que CLIP, proposé par Radford et al. (3), sont entraînés sur de grandes collections
27 de paires image–texte et apprennent à aligner des concepts visuels et linguistiques dans un espace
28 de représentation commun. Grâce à cette supervision sémantique, ces modèles permettent une
29 reconnaissance *zero-shot* et une bonne généralisation à des concepts non vus durant l'entraînement.

30 Bien que les VLM n'aient pas été conçus spécifiquement pour la détection d'objets camouflés, ils
31 offrent une source de connaissances sémantiques complémentaires. Le langage peut fournir des
32 indications de haut niveau sur la présence ou la nature d'un objet, ce qui peut être particulièrement
33 utile lorsque les indices visuels sont ambigus ou peu informatifs.

34 **1.4 Positionnement du projet**

35 Dans ce contexte, ce projet adopte une approche exploratoire visant à combiner trois idées complé-
36 mentaires issues de la littérature : les limites des approches basées sur des images statiques, l'apport
37 du mouvement et de l'information temporelle dans les vidéos, et l'utilisation de connaissances
38 sémantiques fournies par les modèles Vision–Langage. L'objectif du projet n'est pas de proposer
39 une méthode à l'état de l'art, mais plutôt d'étudier comment ces différentes sources d'information
40 peuvent être exploitées conjointement pour améliorer la détection d'animaux camouflés sur le jeu de
41 données MoCA.

42 **References**

- 43 [1] Y. Li, Y. Zhang, Y. Wang, and C. Shen, “Moving camouflaged animals: A dataset and benchmark
44 for motion-based object detection,” in *Proceedings of the IEEE/CVF International Conference*
45 on Computer Vision, 2021. [Online]. Available: <https://arxiv.org/abs/2011.11630v1>
- 46 [2] X. Cheng, J. Fu, F. Zhu, and X. Li, “Implicit motion handling for video camouflaged object
47 detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
48 *Recognition*, 2022. [Online]. Available: <https://arxiv.org/abs/2203.07363>
- 49 [3] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
50 P. Mishkin, and J. Clark, “Learning transferable visual models from natural language supervision,”
51 in *Proceedings of the International Conference on Machine Learning*, 2021. [Online]. Available:
52 <https://arxiv.org/abs/2103.00020>