

# Apprentissage robuste d'une carte de profondeur pour l'évitement d'obstacle dans le cas de caméras volantes, monoculaires et stabilisées

Soutenance de Thèse

---

Clément Pinard

25 juin 2019

U2IS, ENSTA ParisTech

université  
PARIS-SACLAY

ÉCOLE DOCTORALE  
**INTERFACES**  
Approches Interdisciplinaires :  
Fondements, Applications et Innovation



Parrot

1. Introduction et motivations
2. DepthNet Supervisé
3. DepthNet Non Supervisé
4. Utilisation finale de DepthNet
5. Conclusion

# 1. Introduction et motivations

Contexte technologique

Stratégie générale

## 2. DepthNet Supervisé

## 3. DepthNet Non Supervisé

## 4. Utilisation finale de DepthNet

## 5. Conclusion

# Introduction et motivations

---

Contexte technologique

## Le drone grand public

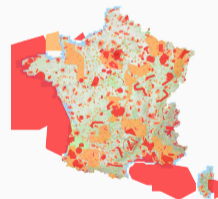
- Accessible et facile à piloter, même pour un néophyte.
- Images aériennes stabilisées et de qualité
- Modes de vol manuels et automatiques

## Accessible à tous ?

- Trop dangereux en cas de crash pour être véritablement grand public !
- La croissance du marché est freinée par des lois de sécurité légitimes mais restrictives



PARROT Anafi (2018)



Carte des zones de vol interdit

## Améliorer la sécurité

- Introduire un évitement d'obstacle automatique
- Utiliser le moins de hardware possible

## Un évitement d'obstacle transparent

- Suivre la consigne, éventuellement directement fixée par l'utilisateur
- Garder la cinématographie en conservant un mouvement fluide



FIGURE 1 – La trajectoire verte est privilégiée

# Introduction et motivations

---

Stratégie générale

## Perception puis évitement

- On se concentre sur la perception des obstacles, avec seulement une caméra monoculaire.
- On essaye de trouver la carte de profondeur  $\theta$  liée à une caméra stabilisée.

## Carte de profondeur

Pour chaque pixel d'une image, donner en plus de la couleur une indication de distance avec la profondeur.



KITTI, mesuré avec un LiDAR [Geiger et al., 2013]

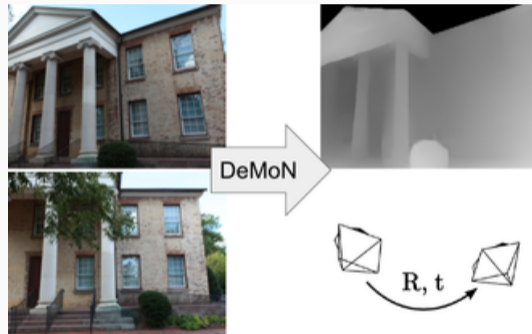


Mesurée avec une Microsoft Kinect [Lai et al., 2011]



Peut-on entraîner à un réseau de neurones à apprendre la profondeur  $\theta$ ?

- Oui!
- Avec plusieurs images : DeMoN [Ummenhofer et al., 2017]
- Avec une seule image : DORN [Fu et al., 2018], MegaDepth [Li and Snavely, 2018]



Fonctionne très bien avec des vidéos de voiture, manque de robustesse pour des scènes inhabituelles.



Ces deux scènes sont plus plates que vous ne le pensez!

Si vous aviez eu une vidéo, vous n'auriez pas eu de problème ...

On ne peut estimer la profondeur  $\tilde{\theta}$  qu'à un **facteur d'échelle** près.

Il faut une grandeur physique disponible.

On propose de définir le facteur d'échelle avec une estimation de la vitesse  $\tilde{v}$ , comparé à la vitesse réelle  $v$ .

	Précédemment [Zhou et al., 2017]	proposition
Predictions	Profondeur $\tilde{\theta}$	Profondeur $\tilde{\theta}$ , Vitesse $\tilde{v}$
Vérité terrain	Profondeur $\theta$	Profondeur $\theta$ , Vitesse $v$
Évaluation	$m = \delta \left( \theta, \tilde{\theta} \times \frac{Me(\theta)}{Me(\tilde{\theta})} \right)$	$m = \delta \left( \theta, \tilde{\theta} \times \frac{ v }{ \tilde{v} } \right)$

**TABLE 1** –  $Me()$  est l'opérateur médian, et  $\delta$  est une fonction de distance (e.g. distance L1)

Cette évaluation permet d'estimer  $\tilde{\theta}$  en supposant la vitesse constante.

1. Introduction et motivations

2. DepthNet Supervisé

Préambule

Estimation du flot optique

StillBox dataset

Résolution du Problème

Resultats

Transfert aux vidéos réelles

3. DepthNet Non Supervisé

4. Utilisation finale de DepthNet

5. Conclusion

# DepthNet Supervisé

---

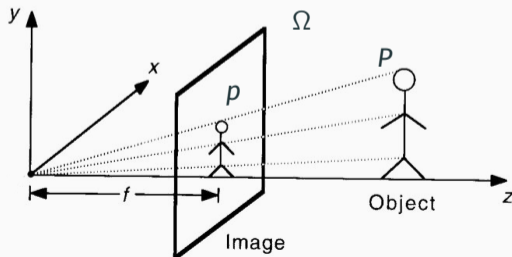
Préambule

## Modèle sténopé

Un point 3D  $P = (X, Y, Z)$  est projeté sur le capteur en un pixel  $p \in \Omega \subset \mathbb{R}^2$ .

Le plan de pixel  $\Omega$  a pour centre optique

$$p_0 = \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}$$



## Mise en équation

$$\begin{aligned} p &= \Pi(KP) \\ \theta(p) &= Z \end{aligned}$$

$$K = \begin{bmatrix} f & 0 & u_0 \\ 0 & f & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{aligned} \Pi : \mathbb{R}^3 &\rightarrow \mathbb{R}^2 \\ P = (x, y, z) &\mapsto \left( \frac{x}{z}, \frac{y}{z} \right) \end{aligned}$$

$$\forall \alpha, \Pi(\alpha P) = \Pi(P)$$

## Passage de $p$ et $\theta$ à $P$

$$P = \theta(p) \times K^{-1} \Pi_{-1}(p)$$

$$\begin{aligned} \Pi_{-1}: \mathbb{R}^2 &\rightarrow \mathbb{R}^3 \\ p = (u, v) &\mapsto (u, v, 1) \end{aligned}$$

## Passage de $P$ à $P_2$

La caméra se déplace avec une rotation  $R$

et une translation  $t = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$

Nouvelle position du point  $P$

$$P_2 = RP + t$$

## Nouvelle position du pixel

$$p_2 = \Pi(KP_2) = \Pi(\theta(p)K RK^{-1} \Pi_{-1}(p) + Kt)$$



## Définition

Le flot optique est un champ de vecteurs représentant le déplacement des pixels d'une image à l'autre.

$$\begin{aligned} F : \Omega &\rightarrow \mathbb{R}^2 \\ p &\mapsto (F_u, F_v) = p_2 - p \end{aligned}$$

- On a déduit le flot optique à partir de la profondeur et du déplacement.
- On cherche à trouver la profondeur  $\theta$  à partir du flot optique et du déplacement.

## Définition

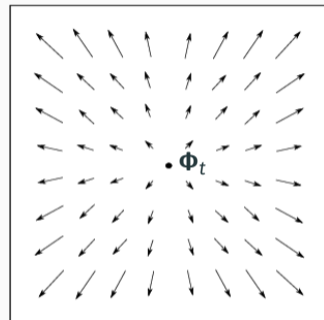
Projection du vecteur translation  $\mathbf{t}$  sur le capteur.

$$\Phi_t = \Pi(\mathbf{Kt})$$

C'est la direction vers laquelle la caméra se dirige

Le flot optique est proportionnel à la profondeur  $\theta$  et à la distance entre  $\Phi_t$  et  $p$

$$\forall p \in \Omega, F(p) = \Pi(\theta(\mathbf{t})\Pi_{-1}(p) + \mathbf{Kt}) = \frac{t_z}{\theta(p) + t_z} (\Phi_t - p)$$



$$\forall p \in \Omega, \theta(p) = \frac{f \|t\|}{\sqrt{f^2 + \|\Phi_t - p_0\|^2}} \left( \frac{\|\Phi_t - p\|}{\|F(p)\|} - \text{sign}(t \cdot u_z) \right)$$

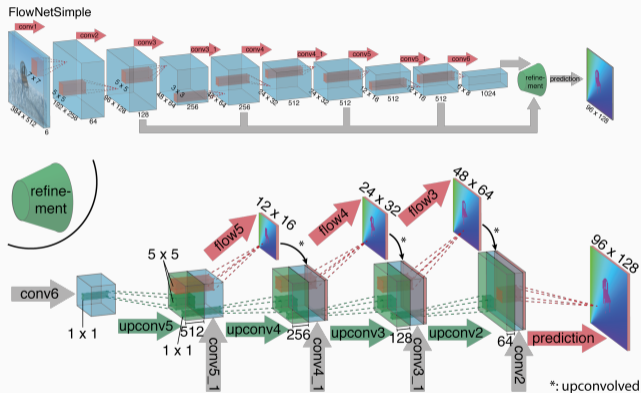
Connaître le flot optique  $F$  et le point  $\Phi_t$  permet de déduire la profondeur.

# DepthNet Supervisé

---

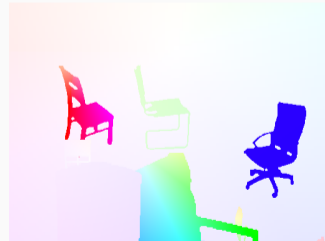
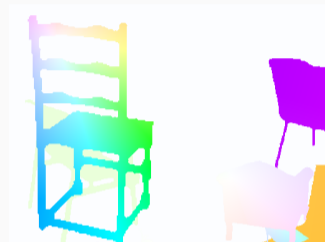
Estimation du flot optique

Réseau Convolutionnel entraîné pour estimer le flot optique d'une paire d'images



Le réseau sort plusieurs échelles du flot optique. Une des architectures les plus reprises car il est très simple.

L'entraînement est fait avec Flying Chairs, un dataset énorme mais irréaliste



$$\forall p \in \Omega, \theta(p) = \frac{f \|t\|}{\sqrt{f^2 + \|\Phi_t - p_0\|^2}} \left( \frac{\|\Phi_t - p\|}{\|F(p)\|} - \text{sign}(t \cdot u_z) \right)$$

Forme indéterminée en  $p = \Phi_t$

Une erreur sur le flot optique estimé  $\tilde{F}$  sera décuplée sur  $\tilde{\theta}$

Le vecteur  $\Phi_t$  doit aussi être correctement estimé

On veut faire apprendre à un réseau de neurones directement la profondeur au lieu d'apprendre le flot optique

# DepthNet Supervisé

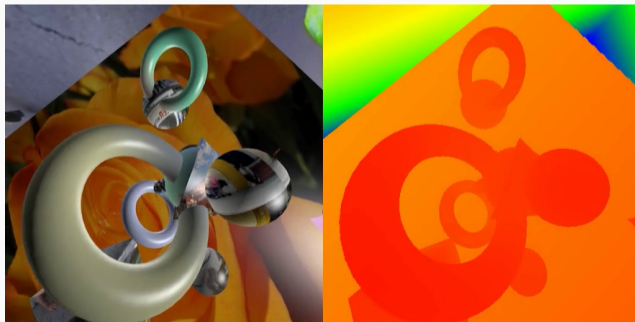
---

StillBox dataset



## Un dataset synthétique représentatif d'une caméra volante stabilisée

- Scènes rigides avec formes primitives en suspension
- Aucune rotation dans le mouvement
- Fait automatiquement sur *Blender*



# DepthNet Supervisé

---

Résolution du Problème

## Architecture

Fully convolutional, très similaire à FlowNet, mais ne sort que la profondeur  $\tilde{\theta}$  à partir de deux images  $I_1$  et  $I_2$

## Apprentissage

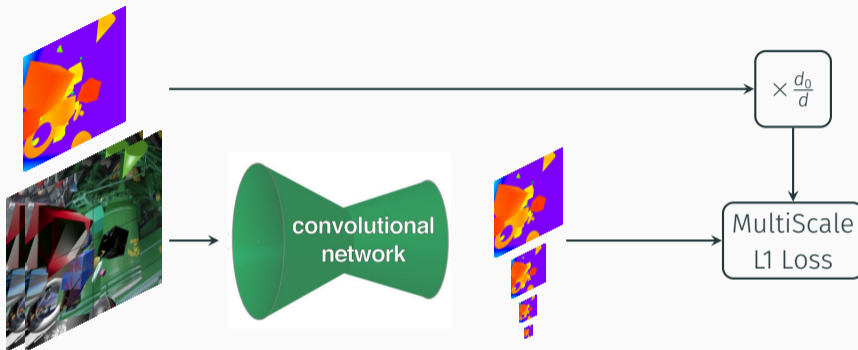
La profondeur  $\tilde{\theta}$  est apprise avec un déplacement nominal  $d_0$  constant.

La profondeur véritable terrain est compensée en fonction du déplacement effectif  $d$  entre  $I_1$  et  $I_2$ . L'estimation  $\tilde{\theta}$  est retrouvée avec l'opération inverse.

$$\theta_{train} = \theta \frac{d_0}{d}$$

$$\tilde{\theta} = \text{DepthNet}(I_1, I_2) \frac{d}{d_0}$$

Le vecteur translation  $\mathbf{t}$  n'est **jamais donné**, seulement sa magnitude  $d$



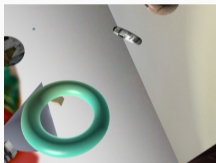
$$\mathcal{L} = \sum_{s \in \text{scales}} \gamma_s \frac{1}{|\Omega|} \sum_{p \in \Omega} |Out_s(p) - \theta_{train}(p)|$$

# DepthNet Supervisé

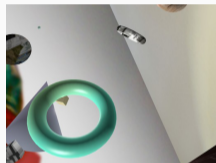
---

Resultats

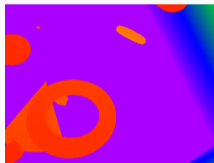
Img 1



Img 2

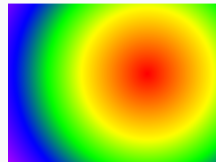
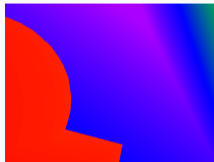
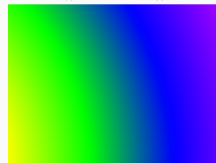


Ground Truth

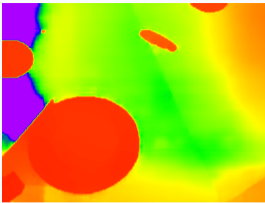


distance to FOE

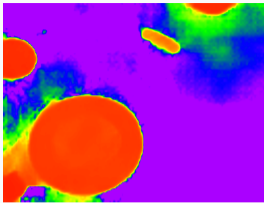
$$\|\Phi_t - p\|$$



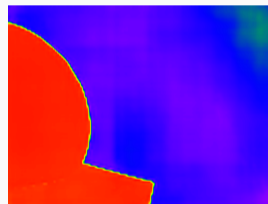
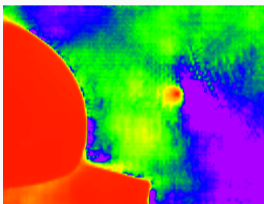
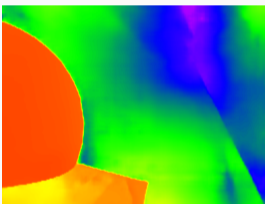
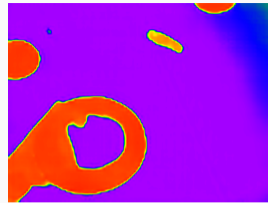
DeMoN



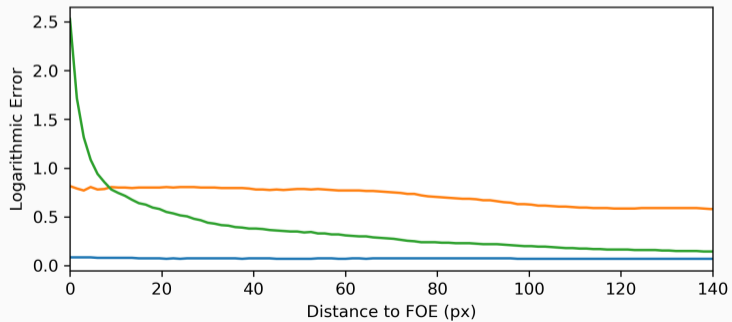
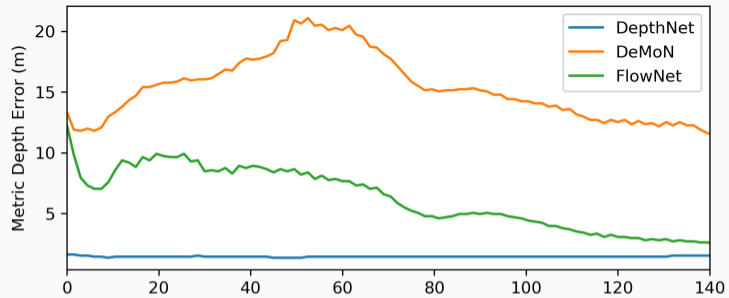
FlowNet



DepthNet

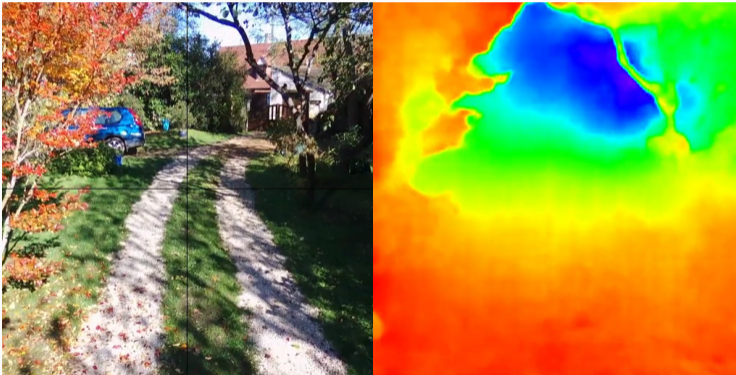


Method	Erreur absolue moyenne (m)	Erreur log moyenne
DepthNet	<b>2.67</b>	<b>0.132</b>
FlowNetS → depth	6.462	0.172
DeMoN	18.85	0.726

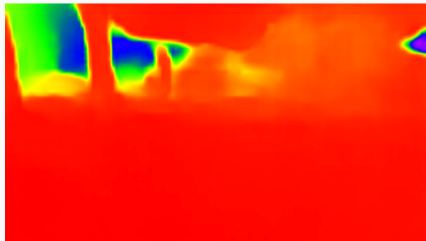
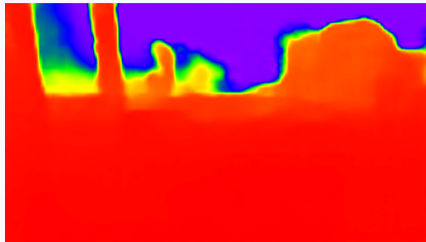




Résultats staisfaisant sur des vidéos stabilisées réelles.



Pas toujours! Problèmes notamment avec le ciel...



Il nous faut un *finetuning* sur des vidéos de drone

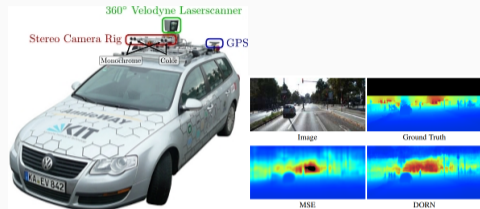
# DepthNet Supervisé

---

Transfert aux vidéos réelles

## Dataset réel

- Construire un dataset avec profondeur est long et compliqué
- La profondeur est souvent incomplète.



## Dataset synthétique

- Plus facile de réaliser beaucoup d'exemples.
- Le rendu peut être photoréaliste [Shah et al., 2017], mais la construction d'une scène complexe est difficile.

L'apprentissage non supervisé permettrait d'avoir un domaine d'apprentissage complexe, réaliste et facile à construire

1. Introduction et motivations

2. DepthNet Supervisé

3. DepthNet Non Supervisé

Préambule

Présentation de SFMLearner

Régularisation par lissage

Régularisation par gestion des occultations

Résultats

4. Utilisation finale de DepthNet

5. Conclusion

# DepthNet Non Supervisé

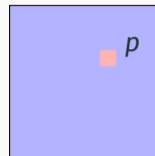
---

Préambule

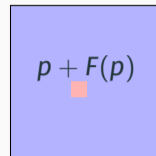
## Hypothèse de la cohérence de la couleur

S'il y a un flot optique  $F$  entre deux images  $I_0$  et  $I_1$ , alors on a une cohérence de la couleur.

$$\forall p \in \Omega, I_0(p) = I_1(p + F(p))$$



$I_0$



$I_1$

## Méthode de Horn et Schunck [Horn and Schunck, 1981]

Tenter de reconstruire  $I_0$  avec les pixels de  $I_1$  pour former  $\tilde{I}_F$ .

$$\forall p \in \Omega, \tilde{I}_F = I_1(p + \tilde{F}(p))$$

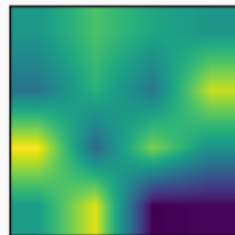
Trouver  $F$  revient à trouver le minimum de la fonction  $\mathcal{L}_p$

$$F = \arg \min_{\tilde{F}} \mathcal{L}_p = \arg \min_{\tilde{F}} (d(I_0, \tilde{I}_{\tilde{F}}))$$

## Optimisation

- On définit l'image sur  $\mathbb{R}^2$  avec une interpolation bilinéaire.
- La construction de  $\tilde{I}_F$  est différentiable en  $\tilde{F}$ .
- On peut minimiser  $\mathcal{L}_p$  par une descente de gradient

$$\tilde{F} \leftarrow \tilde{F} + \mu \nabla_{\tilde{F}} \mathcal{L}_p$$





## Application à la profondeur

On reprend l'équation vue précédemment

$$F = \Pi \left( K R K^{-1} \Pi_{-1}(p) + \frac{Kt}{\theta(p)} \right) - p$$

- Tout est différentiable
- On peut trouver la profondeur  $\theta$  d'une image par optimisation de  $\mathcal{L}_p$ .

# DepthNet Non Supervisé

---

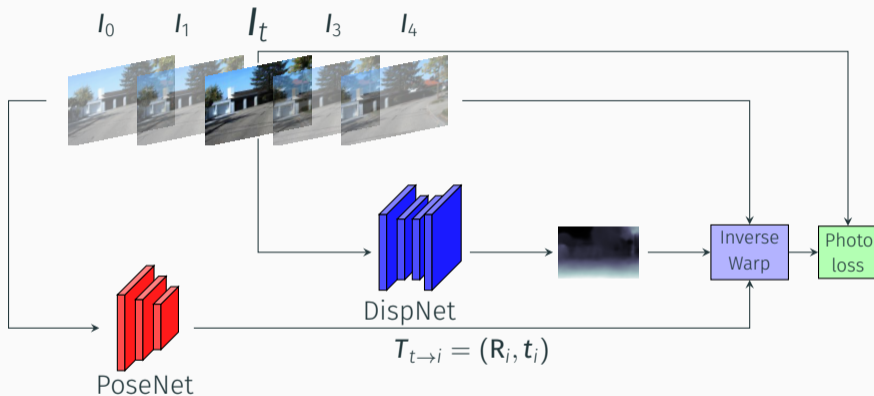
Présentation de SFMLearner

## Solution naïve pour apprendre la profondeur sans la connaître

1. Estimer la profondeur  $\tilde{\theta}$  et le déplacement  $(\mathbf{R}, \mathbf{t})$  par méthodes analytiques de flot optique. Par exemple, par optimisation de la fonction de coût  $\mathcal{L}_p$ .
2. Entraîner DepthNet avec pour consigne  $\theta_{train} = \tilde{\theta} \frac{d_0}{\|\mathbf{t}\|}$

## Idée de SFMLearner

- Entraîner directement un réseau DispNet pour minimiser  $\mathcal{L}_p$
- Entraîner aussi un réseau PoseNet pour calculer le déplacement  $(\mathbf{R}, \mathbf{t})$



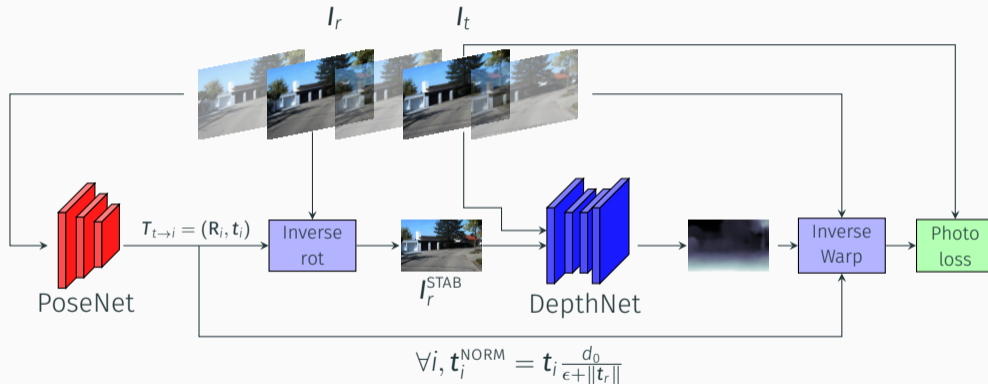
Les réseaux PoseNet et DispNet tentent de minimiser les erreurs photométriques entre  $I_t$  et les images de références déformées  $\tilde{I}_{n \rightarrow t}$

## Stratégie

1. Remplacer DispNet par DepthNet.
2. Choisir dans la séquence une deuxième image  $I_r$  pour former une paire  $(I_r, I_t)$
3. Stabiliser la paire donnée en  $(\tilde{I}_r, I_t)$  pour DepthNet grâce à PoseNet.  
(Pas besoin si les vidéos sont déjà stabilisées)

$$F_{stab} = \Pi (K R K^{-1} \Pi_{-1}(p)) - p$$

4. Normaliser la translation de PoseNet.  $\forall i, \mathbf{t}_i \leftarrow \frac{\mathbf{t}_i d_0}{\|\mathbf{t}_r\|}$

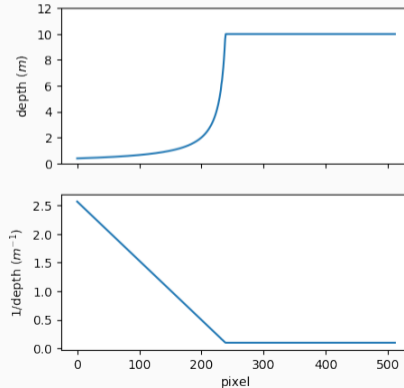
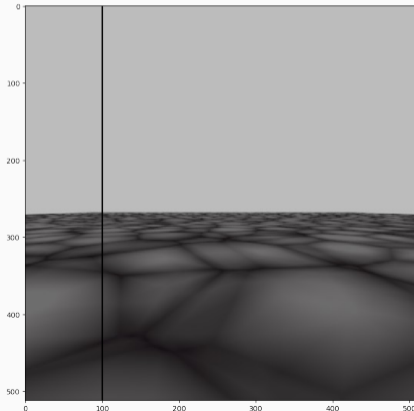


# DepthNet Non Supervisé

---

Régularisation par lissage

Plusieurs solutions possibles permettent d'arriver à une erreur photométrique basse.  
On choisit de contraindre l'optimisation pour avoir la solution la plus physiquement plausible.





## Principe appliqué

Si le gradient d'une image  $\|\nabla I\|$  est bas, alors le Laplacien de sa profondeur inverse  $\Delta\xi = \Delta\left(\frac{1}{\theta}\right)$  doit être bas aussi.

## 2 Idées

1. Rajouter une fonction de coût à l'erreur photométrique pour inciter  $\Delta\xi$  à être bas.
2. Atténuer ce coût en fonction du gradient de l'image  $I$ .

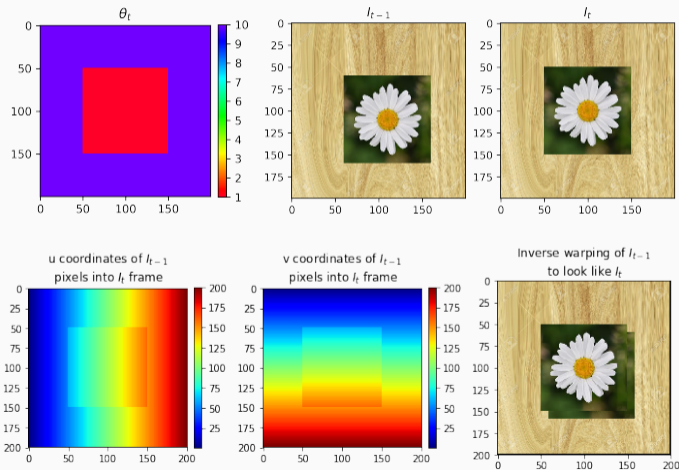
$$\mathcal{L}_{gdiff} = \iint_{\Omega} \exp\left(-\left(\frac{\|\nabla I\|}{\kappa}\right)^2\right) (\Delta\xi)^2 dS$$

## DepthNet Non Supervisé

---

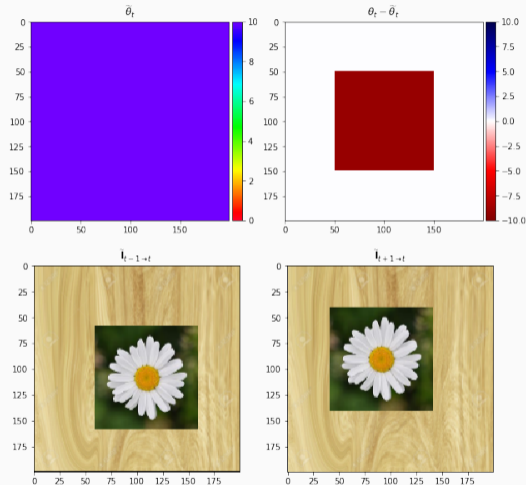
Régularisation par gestion des occultations

Une profondeur parfaite ne minimise pas exactement l'erreur photométrique! Exemple avec une scène typique.



- Ces zones ne doivent pas être optimisées car les pixels correspondant ne sont pas trouvables.
- Un autre, au mauvais endroit mais avec une meilleure couleur sera trouvé.
- Si on peut déduire les zones d'occultation, on les retire de l'erreur photométrique.

$$\mathcal{L}_{pf}(t_2) = \sum_{p \in \Omega} \begin{cases} d(I_t, \tilde{I}_{t_2 \rightarrow t})(p) & \text{si } p \text{ pas} \\ 0 & \text{occulté dans } I_{t_2} \\ & \text{sinon} \end{cases}$$



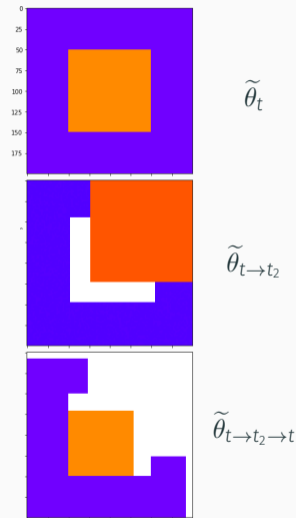
Pour un mouvement particulier :

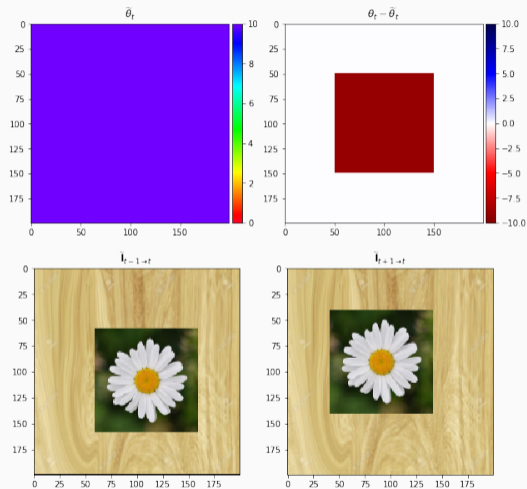
1. Effectuer un warp de la profondeur
2. Effectuer un warp dans l'autre sens

$$\tilde{\theta}_t \xrightarrow{T_{t \rightarrow t_2}} \tilde{\theta}_{t \rightarrow t_2} \xrightarrow{T_{t_2 \rightarrow t}} \tilde{\theta}_{t \rightarrow t_2 \rightarrow t}$$

3. les différences entre  $\tilde{\theta}_t$  et  $\tilde{\theta}_{t \rightarrow t_2 \rightarrow t}$  sont les zones d'occultation.

$$\forall \mathbf{p} \in \Omega, \text{occ}(\mathbf{p}, t_2) = \begin{cases} 1 & \text{si } \tilde{\theta}_{t \rightarrow t_2 \rightarrow t}(\mathbf{p}) \neq \tilde{\theta}_t(\mathbf{p}) \\ 0 & \text{sinon} \end{cases}$$





Erreur photométrique appliquée aux différentes images  $\tilde{I}_{t_i \rightarrow t}$

$$\mathcal{L}_p = \sum_i \sum_{p \in \Omega} d(I_t, \tilde{I}_{t_i \rightarrow t})(p) \times occ(p, t_i)$$

Coût de lissage

$$\mathcal{L}_s = \lambda \mathcal{L}_{gdiff}(\xi, I_t, \kappa) = \frac{\lambda}{|\Omega|} \sum_{p \in \Omega} \exp\left(-\frac{\|\nabla I_t\|^2}{\kappa^2}\right) \times (\Delta \xi)^2$$

Le tout appliqué sur les différentes échelles de sortie  $\tilde{\theta}^s$  du réseau

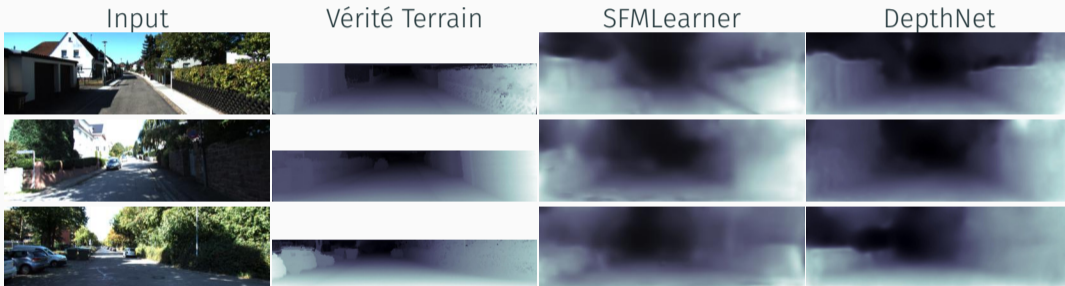
$$\mathcal{L} = \sum_s \frac{1}{2^s} (\mathcal{L}_p^s + \mathcal{L}_s^s)$$

# DepthNet Non Supervisé

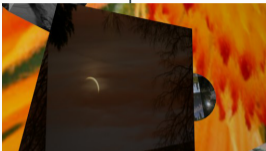
---

Résultats





Input



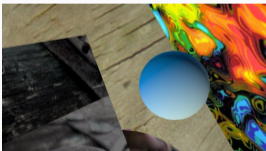
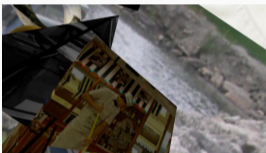
Vérité Terrain

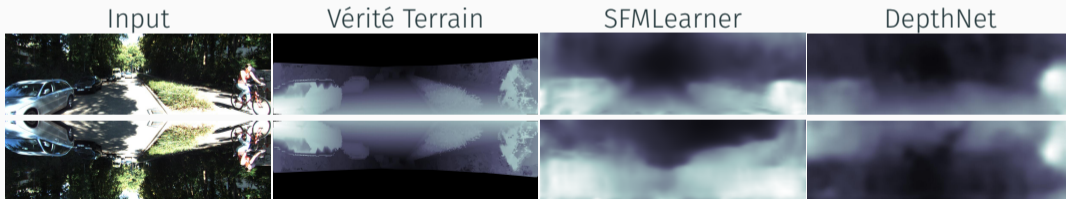


SFMLearner



DepthNet





## Fonctions d'erreur

Erreur Absolue  $|\theta - \tilde{\theta}|$

Erreur Relative  $\frac{|\theta - \tilde{\theta}|}{\theta}$

Erreur Logarithmique  $\left| \ln \left( \frac{\theta}{\tilde{\theta}} \right) \right|$

L'évaluation est fait en prenant la moyenne de ces fonctions sur un ensemble de validation.

## Fonctions de précision

Proportion de valeurs dans un intervalle de confiance autour de  $\theta$ .

$$P_{\delta} = P \left( \tilde{\theta} \in \left[ \frac{1}{\delta} \theta, \delta \theta \right] \right)$$

$P_{1.25}$  proportion de points avec moins de **25%** d'erreur

$P_{1.25^2} = P_{1.56}$  moins de **56%** d'erreur

$P_{1.25^3} = P_{1.95}$  moins de **95%** d'erreur

Method	scale factor	Average error	Relative error	Log error	$P_{1.25}$	$P_{1.25^2}$	$P_{1.25^3}$
SFMLearner	GT	3.455	0.185	0.191	0.714	0.899	0.962
SFMLearner	P	4.629	0.296	0.287	0.542	0.817	0.918
DepthNet 1	P	7.872	0.650	0.446	0.356	0.613	0.775
DepthNet 2	P	4.078	0.268	0.222	0.683	<b>0.883</b>	<b>0.944</b>
DepthNet 3	P	<b>4.001</b>	<b>0.262</b>	<b>0.221</b>	<b>0.687</b>	<b>0.883</b>	0.943

## Rappel

 doit être bas

 doit être haut

DepthNet 1 : entraînement supervisé sur dataset synthétique (StillBox)

DepthNet 2 : entraînement non supervisé sur KITTI

DepthNet 3 : entraînement non supervisé sur KITTI, images stabilisées

Method	scale factor	Detecteur d'occultation	Average error	Relative error	Log error	$P_{1.25}$	$P_{1.25^2}$	$P_{1.25^3}$
DepthNet 1	P	✗	<b>3.822</b>	0.224	0.229	0.718	0.856	0.927
SFMLearner	GT	✗	11.43	0.593	0.633	0.479	0.656	0.756
SFMLearner	P	✗	14.32	0.816	0.461	0.302	0.542	0.696
DepthNet 2	P	✗	5.248	0.259	0.232	0.709	0.887	0.928
DepthNet 2	P	✓	5.060	<b>0.212</b>	<b>0.206</b>	<b>0.729</b>	<b>0.899</b>	<b>0.943</b>

DepthNet 1 : entraînement supervisé

DepthNet 2 : entraînement non supervisé



0m

10m

150m

1. Introduction et motivations

2. DepthNet Supervisé

3. DepthNet Non Supervisé

4. Utilisation finale de DepthNet

Déplacement optimal

Multi-inférence

Preuve de concept avec streaming vidéo sur Bebop

5. Conclusion



## Utilisation finale de DepthNet

---

Déplacement optimal

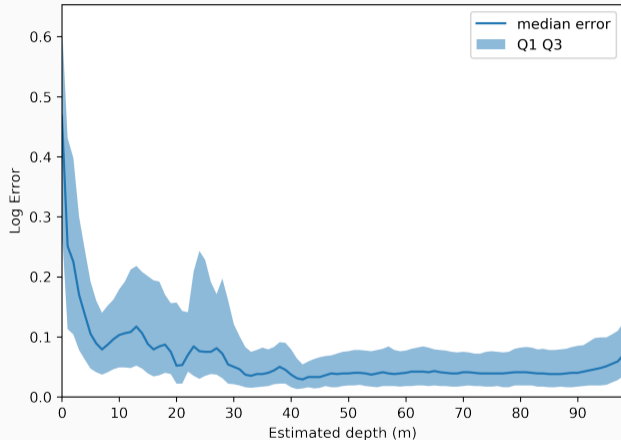
## De DepthNet à la profondeur $\tilde{\theta}$

On multiplie par le ratio en  $d_0$  et le déplacement mesuré par le drone  $d(t, \delta t)$

$$\tilde{\theta}(t) = \text{DepthNet}(I_t, I_{t-\delta t}) \frac{d(t, \delta t)}{d_0}$$

$$d(t, \delta t) = \left\| \int_{t-\delta t}^t \mathbf{V}(\tau) d\tau \right\|$$

- La sortie de DepthNet est inversement proportionnelle au déplacement  $d$ .
- On peut jouer sur le décalage temporel  $\delta_t$  pour contrôler la distribution de DepthNet.
- On a 30 images par seconde possibles pour jouer sur le paramètre  $\delta t$



Toutes les valeurs de sortie de DepthNet ne se valent pas!  
On veut que Depthnet soit à valeurs autour de 70m pour qu'il fonctionne sous les meilleures conditions possibles

On introduit la fonction  $\beta$  à valeurs dans  $[0, 1]$  et les paramètres  $\alpha$  et  $\theta_{max}$  pour simplifier cette expression.

$$\begin{aligned} \beta : \mathbb{R}^{6|\Omega|} &\rightarrow [0, 1]^{|\Omega|} & \alpha &= \frac{\theta_{max}}{d_0} & \theta_{max} &= 100 \\ I_1, I_2 &\mapsto \frac{1}{\theta_{max}} \text{DepthNet}(I_1, I_2) \end{aligned}$$

$$\tilde{\theta}(t) = \alpha\beta(I_t, I_{t-\delta t})d(t, \delta t) = \alpha\beta d(t, \delta t)$$

On veut que  $\beta$  ait une moyenne de  $\beta_{mean} = 0.7$

$$d_{opt} = d_0 \frac{\mathbb{E}(\beta)}{\beta_{mean}} = \frac{\mathbb{E}(\tilde{\theta})}{\alpha\beta_{mean}}$$

## Stratégie

1. Estimer  $\tilde{\theta}$  avec une certaine paire d'images  $I_t, I_{t-\delta_t}$
2. Calculer la moyenne de  $\beta$
3. En déduire le déplacement optimal  $d_{opt}$
4. Trouver le nouveau décalage  $\delta'_t$  avec un déplacement suffisamment proche de  $d_{opt}$

## Utilisation finale de DepthNet

---

Multi-inférence



Certaines scènes n'ont pas qu'un plan  
Il n'y a pas un décalage optimal mais **plusieurs**.

## Principe de la multi-inférence

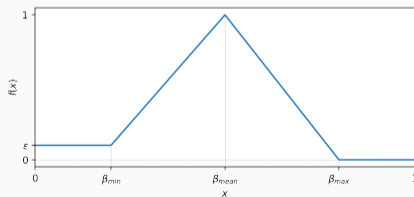
- Trouver  $n$  centroïdes  $c_0, \dots, c_{n-1}$  à partir de l'estimation  $\tilde{\theta}$  précédente et un algorithme type K-means [MacQueen, 1967]
- En déduire les décalages optimaux  $\delta t_0, \dots, \delta t_{n-1}$
- Fusionner les différentes estimations  
 $\beta_0, \dots, \beta_{n-1}$

## Stratégie

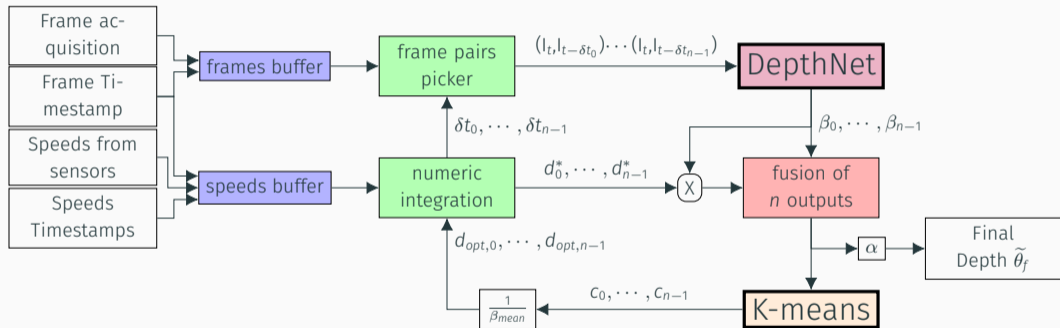
- On compense les  $\beta_i$  avec le déplacement correspondant
- On fait une moyenne pondérée pixel par pixel des  $\tilde{\theta}_i$
- La pondération est forte quand la valeur  $\beta_i(\mathbf{p})$  est proche de  $\beta_{mean}$ , faible autrement.

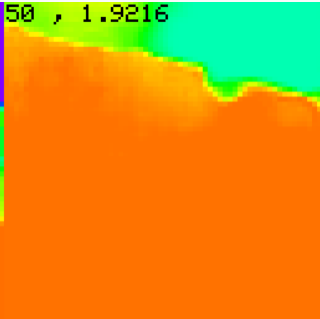
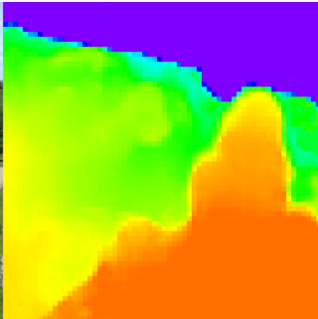
$$\forall \mathbf{p} \in \Omega, \tilde{\theta}_f(\mathbf{p}) = \alpha \frac{\sum_i d_i^* w_i(\mathbf{p}) \beta_i(\mathbf{p})}{\sum_i w_i(\mathbf{p})}$$

$$w_i(\mathbf{p}) = f(\beta_i(\mathbf{p}))$$

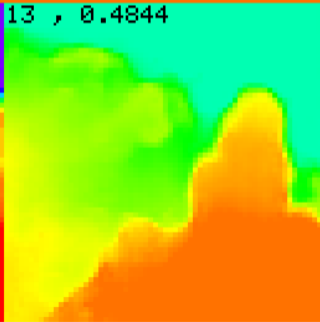
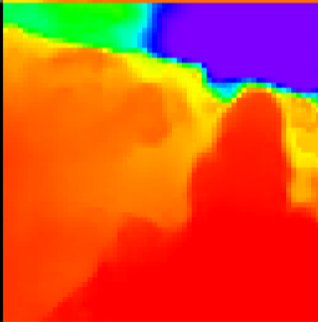
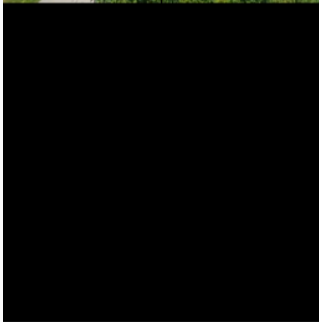








50 , 1.9216

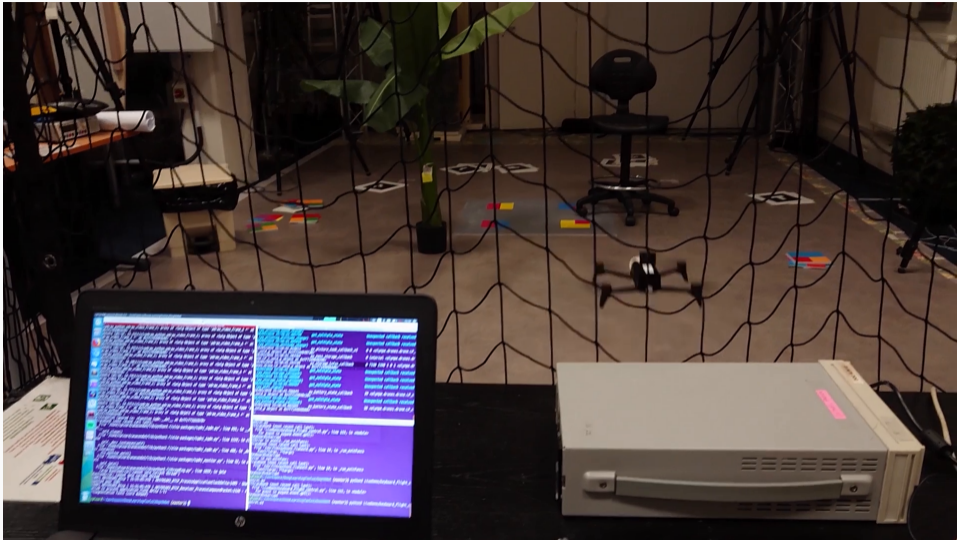


13 , 0.4844

## Utilisation finale de DepthNet

---

Preuve de concept avec streaming vidéo  
sur Bebop



1. Introduction et motivations

2. DepthNet Supervisé

3. DepthNet Non Supervisé

4. Utilisation finale de DepthNet

5. Conclusion

- Introduction d'un réseau de neurones pour l'estimation de profondeur. Réseau axé sur la robustesse aux scènes inhabituelles. (*Article et présentation orale à UAVg 2017, Brevet pour Parrot déposé en 2017*)
- Création de Still Box, un dataset synthétique où la robustesse est essentielle. (*Article et présentation orale à UAVg 2017*)
- Adaptation d'un algorithme connu SFMLearner pour apprendre DepthNet sans aucune supervision. (*Workshop à ECCV 2018*)
- Introduction d'un algorithme pour utiliser le réseau en conditions réelles et avoir une portée arbitrairement grande. (*Article à ECMR 2017*)

- Création d'un dataset de validation de profondeur pour les vidéos de drone.
- Gestion des scènes mobiles.
- Optimisation de l'architecture de DepthNet pour l'embarqué




Merci!




Code source disponible sur github




`perso.ensta.fr/~pinard/`





-  Dosovitskiy, A., Fischer, P., Ilg, E., Häusser, P., Hazırbaş, C., Golkov, V., v.d. Smagt, P., Cremers, D., and Brox, T. (2015).  
**FlowNet : Learning optical flow with convolutional networks.**  
*In IEEE International Conference on Computer Vision (ICCV).*
-  Fu, H., Gong, M., Wang, C., Batmanghelich, K., and Tao, D. (2018).  
**Deep ordinal regression network for monocular depth estimation.**  
*In IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*
-  Geiger, A., Lenz, P., Stiller, C., and Urtasun, R. (2013).  
**Vision meets robotics : The kitti dataset.**  
*The International Journal of Robotics Research*, 32(11) :1231–1237.

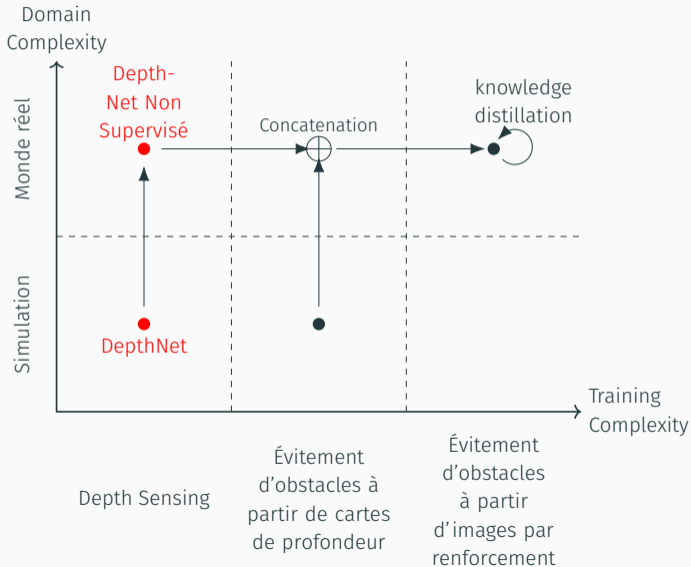
-  Horn, B. K. and Schunck, B. G. (1981).  
**Determining optical flow.**  
*Artificial intelligence*, 17(1-3) :185–203.
-  Lai, K., Bo, L., Ren, X., and Fox, D. (2011).  
**A large-scale hierarchical multi-view rgb-d object dataset.**  
In *2011 IEEE international conference on robotics and automation*, pages 1817–1824.  
IEEE.
-  Li, Z. and Snavely, N. (2018).  
**Megadepth : Learning single-view depth prediction from internet photos.**  
In *Computer Vision and Pattern Recognition (CVPR)*.

-  MacQueen, J. (1967).  
**Some methods for classification and analysis of multivariate observations.**  
In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1 : Statistics*, pages 281–297, Berkeley, Calif. University of California Press.
-  Shah, S., Dey, D., Lovett, C., and Kapoor, A. (2017).  
**AirSim : High-fidelity visual and physical simulation for autonomous vehicles.**  
In *Field and Service Robotics*.
-  Ummenhofer, B., Zhou, H., Uhrig, J., Mayer, N., Ilg, E., Dosovitskiy, A., and Brox, T. (2017).  
**Demon : Depth and motion network for learning monocular stereo.**  
In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.



- Zhou, T., Brown, M., Snavely, N., and Lowe, D. G. (2017).  
**Unsupervised learning of depth and ego-motion from video.**  
In *CVPR*.

# Résumé de la stratégie



# Apprentissage profond

## Principe général

Fonction qui dépend d'une entrée (par exemple une image  $I$ ) et de paramètres internes appelés poids ( $W$ )

$$\text{output} = f(W, I)$$

## Apprentissage

Minimiser une fonction de coût  $\mathcal{L}$ .

Dans le cas d'un apprentissage, on minimise la différence avec une consigne  $t$ .

$$W_{\text{final}} = \arg \min_W \mathbb{E}_{(I,t) \sim \mathcal{V}} \mathcal{L}(f(W, I), t)$$

Pour y arriver on fait un descente de gradient

$$W \leftarrow W - \lambda \nabla_W \mathcal{L}(f(W, I), t)$$

$$\mathcal{L}_{diff}(\xi) = \iint_{\Omega} \exp\left(-\left(\frac{\|\nabla I\|}{\kappa}\right)^2\right) \|\nabla \xi\|^2 dS \quad (1)$$

$$\mathcal{L}_{TV}(\xi) = \iint_{\Omega} \exp\left(-\frac{\|\nabla I\|}{\kappa}\right) \|\nabla \xi\| dS \quad (2)$$

$$\mathcal{L}_{gdiff}(\nabla \xi) = \iint_{\Omega} \exp\left(-\left(\frac{\|\nabla I\|}{\kappa}\right)^2\right) (\Delta \xi)^2 dS \quad (3)$$

$$\mathcal{L}_{TVV}(\nabla \xi) = \iint_{\Omega} \exp\left(-\frac{\|\nabla I\|}{\kappa}\right) |\Delta \xi| dS \quad (4)$$

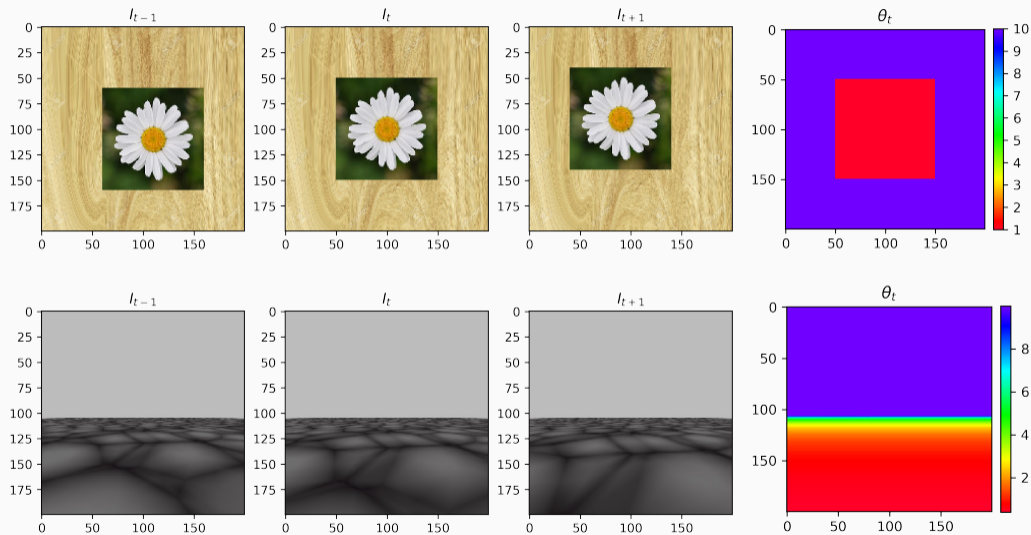
$$\mathcal{L}(\tilde{\theta}, \alpha) = \mathcal{L}_p(\tilde{\theta}) + \lambda \mathcal{L}_s(\alpha) + \gamma d\left(\alpha, \xi = \frac{1}{\tilde{\theta}}\right)$$

## Stratégie d'optimisation "robuste"

1. estimer  $\tilde{\theta} = f(W)$  avec les paramètres optimisables  $W$ .
2. optimiser  $\mathcal{L}(\tilde{\theta}, \alpha)$  avec  $\tilde{\theta}$  fixé, trouver  $\tilde{\alpha} = \arg \min_{\alpha} \mathcal{L}(\tilde{\theta}, \alpha)$ .
3. Caculer  $\nabla_W \mathcal{L}(\tilde{\theta}, \tilde{\alpha})$  et mettre à jour  $W$ .



# Scènes synthétiques



$$LE(\theta, \tilde{\theta}) = \left\| \log(\theta) - \log(\tilde{\theta}) \right\|$$

$$LE_0 = LE(\tilde{\theta}_0^S, \theta^S)$$

$$\mathcal{M}_S(\tilde{\theta}^S) = \begin{cases} 0 & \text{if } LE(\theta^S, \tilde{\theta}^S) > LE_0 \\ 1 - \frac{LE(\theta^S, \tilde{\theta}^S)}{LE_0} & \text{otherwise} \end{cases}$$

$$\mathcal{M} = \sqrt[|S|]{\prod_{S \in \mathcal{S}} \mathcal{M}_S}$$

## Rechercher d'hyperparamètres sur scènes synthétiques

- Maximiser  $\mathcal{M}$
- Minimiser  $\Delta\mathcal{M}$  : le maximum de  $\mathcal{M}$  doit être le plus étendu possible

# Resultats

