

PSRA - raport

Algorytm Flux Tensor with Split Gaussian models

Piotr Janus, Kamil Piszczek

1. Wprowadzenie

Opracowanie i implementacja sprzętowa algorytmu *FTSG*, było głównym celem przedstawianego projektu. Algorytm ten jest metodą hybrydową, dzięki czemu posiada wszystkie zalety algorytmów *Flux Tensor* (rozdział 2) oraz *Split Gaussian Models* (rozdział ??).

Obie wykorzystywane metody nie są idealne, jednak wzajemnie się uzupełniają przez co można opracować algorytm wykorzystujący je jednocześnie. Metoda *Flux Tensor* nie daje możliwości wykrywania obiektów statycznych. Zapewnia jednak odporność na zakłócenia spowodowane zmianami oświetlenia, a częściowo także przez występowanie cieni. Przeciwnieństwem jest algorytm wykorzystujący rozkłady Gaussa, o ograniczonej odporności na zakłócenia, ale dający możliwość wykrywania obiektów statycznych.

2. Metoda Flux Tensor

2.1. Opis Algorytmu

Algorytm *Flux Tensor* jest uniwersalną metodą wykrywania ruchomych obiektów na obrazie ze statycznej kamery. Metoda ta jest niewrażliwa na zmiany oświetlenia, jednak przy jej użyciu nie jesteśmy w stanie rozpoznać statycznych obiektów pierwszoplanowych. Algorytm opiera się na detekcji krawędzi poziomych i pionowych, a następnie obliczeniu pochodnej obrazu po czasie, w celu wykrycia ruchomych elementów.

Metoda *Flux Tensor* w przedstawionej wersji operuje na obrazie w skali szarości. Oznaczmy przez $I(x, y, t)$ wartość piksela wejściowego. Gdzie x, y to współrzędne piksela na obrazie, t – czas. Dla uproszczenia wektor w przestrzeni (x, y, t) możemy zapisać jako v . *Flux Tensor* możemy przedstawić, jako macierz, opisującą zmiany wartości piksela wejściowego w lokalnej trójwymiarowej czasoprzestrzeni. Przez Ω rozumiemy otoczenie w przestrzeni (x, y, t) , które rozważamy przy analizie danego piksela. Zgodnie z artykułem [?] macierz Flux Tensor wygląda następująco:

$$J_F(i) = \begin{pmatrix} \int_{\Omega} \left(\frac{\partial^2 I(v)}{\partial x \partial t} \right)^2 dv & \int_{\Omega} \frac{\partial^2 I(v)}{\partial x \partial t} \frac{\partial^2 I(v)}{\partial y \partial t} dv & \int_{\Omega} \frac{\partial^2 I(v)}{\partial x \partial t} \frac{\partial^2 I(v)}{\partial^2 t} dv \\ \int_{\Omega} \frac{\partial^2 I(v)}{\partial y \partial t} \frac{\partial^2 I(v)}{\partial x \partial t} dv & \int_{\Omega} \left(\frac{\partial^2 I(v)}{\partial y \partial t} \right)^2 dv & \int_{\Omega} \frac{\partial^2 I(v)}{\partial y \partial t} \frac{\partial^2 I(v)}{\partial^2 t} dv \\ \int_{\Omega} \frac{\partial^2 I(v)}{\partial^2 t} \frac{\partial^2 I(v)}{\partial x \partial t} dv & \int_{\Omega} \frac{\partial^2 I(v)}{\partial^2 t} \frac{\partial^2 I(v)}{\partial y \partial t} dv & \int_{\Omega} \left(\frac{\partial^2 I(v)}{\partial^2 t} \right)^2 dv \end{pmatrix} \quad (1)$$

Elementy macierzy zawierają informację na temat zmiany gradientu w czasie, co pozwala na rozróżnienie elementów ruchomych i statycznych na obrazie. Detekcja odbywa się bezpośrednio na podstawie śladu macierzy:

$$trace(J_F) = \int_{\Omega} \left\| \frac{\partial}{\partial t} \nabla I(v) \right\|^2 dv \quad (2)$$

Ślad macierzy jest obliczany dla każdego piksela. Jeżeli jego wartość jest większa od ustalonej wartości progowej T , to analizowany piksel jest elementem ruchomego obiektu. W przeciwnym wypadku zakładamy, że piksel jest tłem. Przyjmijmy następujące oznaczenia:

$$I_{xt} = \frac{\partial^2 I(x, y, t)}{\partial x \partial t}, \quad I_{yt} = \frac{\partial^2 I(x, y, t)}{\partial y \partial t}, \quad I_{tt} = \frac{\partial^2 I(x, y, t)}{\partial^2 t} \quad (3)$$

Równanie (2) można wówczas zapisać jako:

$$trace(J_F) = \int_{\Omega} (I_{xt}^2(i) + I_{yt}^2(i) + I_{tt}^2(i)) dv \quad (4)$$

Składowe I_{xt} i I_{yt} zawierają informację o wykrytych i jednocześnie poruszających się krawędziach (na podstawie pochodnej po czasie), natomiast I_{tt} informuje o zmianach oświetlenia w czasie.

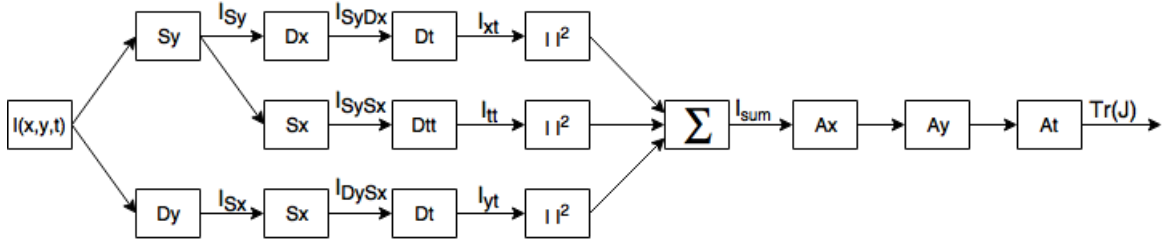
2.2. Obliczenia numeryczne

W pierwszym kroku piksel w formacie *RGB* konwertowany jest do skali szarości zgodnie z zależnością:

$$GRAY = 0.21R + 0.72G + 0.07B \quad (5)$$

Składowe I_{xt} , I_{yt} , I_{tt} są obliczane z wykorzystaniem jednowymiarowych masek konwolucji. Jest to rozwiązanie mniej złożone obliczeniowo, niż stosowanie masek 3D. Dodatkowo, w celu redukcji szumów, używany jest jednowymiarowy filtr wygładzający dla składowej, która nie występuje w aktualnie obliczanej pochodnej.

Zgodnie z powyższym, podczas obliczania składowej I_{xt} są używane filtry różniczkujące dla składowych x i t , oraz filtr wygładzający dla składowej y . Analogicznie, w przypadku składowej I_{yt} należy wykorzystać filtr wygładzający dla składowej x . Ostatnim elementem jest składowa I_{tt} . Tutaj filtry wygładzające używane są dla składowych x i y . Operacja całkowania jest wykonywana z użyciem odpowiednich jednowymiarowych filtrów. Na obraz nakładane są maski uśredniające dla poszczególnych składowych. Schemat tych operacji został przedstawiony na rys. 1.



Rysunek 1: Schemat operacji wyznaczających ślad macierzy J_F

Przedstawione na rys. 1 operacje to odpowiednio:

- Dx , Dy – filtry różniczkujące odpowiednio dla składowych x i y
- Sx , Sy – filtry wygładzające dla składowych x i y
- Dt , Dtt – odpowiednio pierwsza i druga pochodna z wartości piksela po czasie
- Ax , Ay – filtry uśredniające przestrzenne
- At – filtr uśredniający czasowy

Oprócz rozmiarów poszczególnych filtrów algorytm posiada także parametr T , czyli próg wartości śladu, macierzy, powyżej którego piksel zostaje uznany za tło. Ostateczna lista parametrów algorytmu wygląda następująco:

- T –wartość progowa
- nDs – rozmiar filtrów przestrzennych (różniczkujących i wygładzających)
- nDt – liczba ramek używana do różniczkowania po czasie
- nAs – rozmiar masek uśredniających przestrzennych
- nAt – liczba ramek używana do uśredniania po czasie

Należy zwrócić uwagę, że dobranie dużych rozmiarów filtrów (szczególnie różniczkujących po czasie) znacząco zwiększa zużycie pamięci. Domyślnymi parametrami, sprawdzającymi się w większości przypadków są:

$$T = 25, \quad nDs = 5, \quad nDt = 5, \quad nAs = 5, \quad nAt = 5$$

Jak wcześniej wspomniano, schemat na rys. 1, zoptymalizowany został pod kątem zużycia pamięci z zachowaniem równoległości obliczeń. W tym przypadku do działania algorytmu będą potrzebne 3 bufor o rozmiarze $nDt - 1$ każdy oraz jeden bufor o rozmiarze $nAt - 1$. Bufory działają na zasadzie kolejki *FIFO* (ang. *First In, First Out* – pierwszy na wejściu, pierwszy na wyjściu). W przypadku domyślnych parametrów wiąże się to z przechowywaniem w pamięci 16 ramek obrazu.

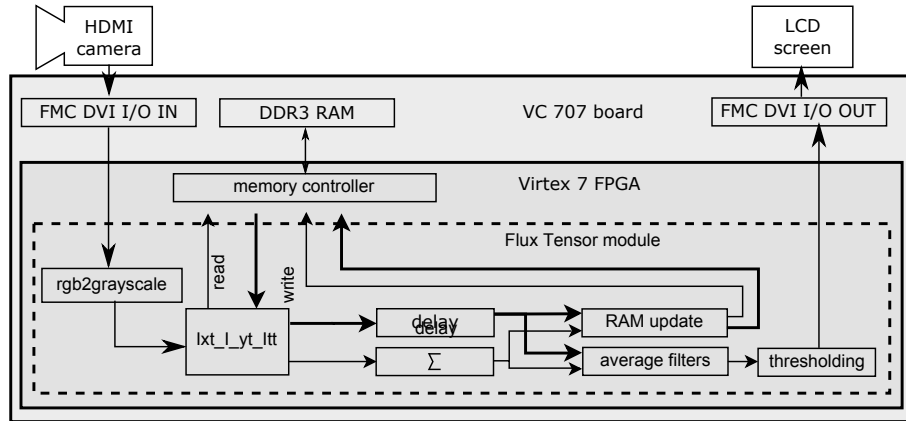
2.3. Implementacja sprzętowa

Wstępna analiza algorytmu przedstawionego na rysunku 1 pozwoliła wskazać dwa podstawowe wyzwania związane z realizacją sprzętową. Po pierwsze, z uwagi na konieczność buforowania 16 ramek obrazu (domyślne parametry) niezbędne okazało się być użycie zewnętrznej pamięci RAM – wspólnie dostępne układy FPGA nie posiadają wystarczających zasobów pamięciowych do buforowania ramek HD. Po drugie algorytm wymaga wykonania wielu operacji kontekstowych. W celu zaoszczędzenia zasobów, warto ustawić operacje w takiej kolejności, aby pierwsze wykonywane były operacje wykorzystujące kontekst pionowy.

Jak zaznaczono wcześniej, wyliczanie różniczki i uśrednianie po czasie (operacje D_t, D_{tt}, A_t) wymagają buforowania poprzednich ramek – wartości $I_{S_y D_x}, I_{S_y S_x}, I_{D_y S_x}$ oraz I_{sum} . Liczba zapamiętywanych wartości zależy od zastosowanych rozmiarów masek. Dla podanych wcześniej parametrów (wszystkie maski o rozmiarze 5), dla pojedynczego piksela wymaga się $(5 - 1) * 4 * 8$ bitów = 128 bitów. Dla maksymalnego rozważanego rozmiaru maski $n = 7$ wskaźnik ten wynosi 192 bity. Wyznaczone wartości nie przekraczają możliwości współczesnych układów FPGA współpracujących z szybkimi pamięciami DDR3 lub DDR4, nawet dla rozdzielczości HD.

2.3.1. Stworzona architektura

Ogólny schemat zrealizowanego modułu sprzętowego przedstawiono na rysunku 2. Wszystkie moduły zostały opisane w języku Verilog, w trakcie prac wykorzystywano środowisko Vivado w wersji 2015.4 firmy Xilinx.



Rysunek 2: Architektura zaimplementowana w układzie FPGA

Źródłem sygnału wizyjnego była kamera Sony HDR-CX280 z wyjściem HDMI. Do wizualizacji rezultatów działania systemu wykorzystano typowy monitor LCD z wejściem HDMI. Wszystkie obliczenia realizowane byłyby na karcie ewaluacyjnej VC 707 firmy Xilinx. Centralnym elementem platformy jest układ FPGA Virtex 7 (XC7VX485T-2FFG1761). W rozwiązaniu wykorzystano ponadto zewnętrzną pamięć DDR3 RAM oraz moduł wejścia/wyjścia HDMI (przystawka Avnet FMC DVI I/O).

W pierwszym kroku sygnał wejściowy konwertowany jest z przestrzeni barw RGB do skali szarości (moduł **rgb2grayscale**). Zastosowano przekształcenie zgodne z biblioteką OpenCV.

W następnym module wyznaczone są składowe I_{xt}, I_{yt}, I_{tt} (szczegółowy opis modułu w Subsection 2.3.2). Obliczanie pochodnych po czasie wymaga wykorzystania kontrolera pamięci (moduł **memory controller**). Następnie wyliczane wartości są podnoszone do kwadratu oraz sumowane (moduł Σ). Równolegle następuje opóźnienie danych odczytanych z pamięci RAM.

Wyniki poprzednich operacji trafiają do dwóch modułów:

- **RAM update** – następuje tutaj aktualizacja danych, które mają zostać zapisane do pamięci RAM. Polega ona na operacji FIFO (First In First Out) tj. usunięciu najstarszego elementu z wektora danych,
- **average filters** – realizacja uśredniania (filtry A_x, A_y, A_t) – omówienie w Subsection 2.3.3.

Ostatnim etapem algorytmu jest binaryzacja (moduł **thresholding**).

W trakcie projektowania modułu dążono do redukcji zasobów poprzez ograniczenie długości linii opóźniających dla danych zapisanych w pamięci RAM (słowo o szerokości 128 bitów). Po pierwsze dane czytano dokładnie w momencie, kiedy były potrzebne. Ponadto zdecydowano się na zapisanie do pamięci wartości I_{sum} , nie I_{sum} po operacjach A_x i A_y . Zastosowanie takiego podejścia wiąże się z koniecznością realizacji filtracji A_x i A_y dla wszystkich pikseli z kontekstu czasowego (tj. 5 w rozważanym przypadku). Jeśli weźmie się pod uwagę tylko długie linie opóźniające (zrealizowane w pamięci BRAM), to przy takim rozwiązaniu niezbędne jest pamiętanie: $YY \cdot 4 \cdot 5 \cdot 8 = 160 \cdot XX$ bitów (YY – rozdzielczość Alternatywne rozwiązanie, tj. bezpośrednia realizacja wg. schematu z rysunku 1, wymaga pamiętania $YY \cdot 4 \cdot 128 = 512 \cdot YY$ bitów).

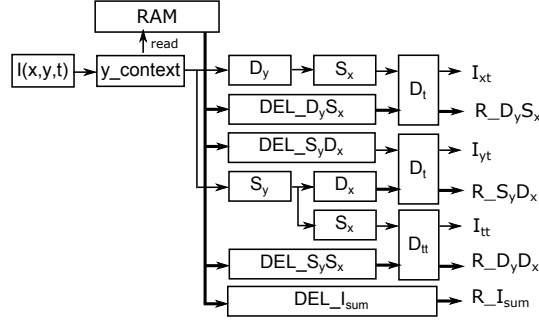
Warto też zaznaczyć, że zastosowane skrócenie ścieżki może mieć również inne zalety. Przykładowo, jeśli moduł FT wykorzystuje się w ramach bardziej rozbudowanego systemu wizyjnego, to brak konieczności opóźniania pozostałych danych zapisanych w pamięci RAM jest bardzo pożądanym i może pozwolić na dużo większą redukcję zapotrzebowania na zasoby.

2.3.2. Moduł obliczający ślad macierzy

Szczegółowy schemat modułu $I_{xt}I_{ty}I_{tt}$ przedstawiono na rysunku 3. Na początku wyliczany jest kontekst pionowy. W momencie gdy na wyjściu modułu pojawi się poprawny kontekst tj. wektor 5 pikseli, generowany jest sygnał odczytu z pamięci RAM (**read**), tak aby dostępne były zapisane dla rozważanego piksela dane. Mają one następującą postać:

$$[I_{S_y D_x}(t-1), \dots, I_{S_y D_x}(t-nDt+1), I_{D_y S_x}(t-1), \dots, I_{D_y S_x}(t-nDt+1), \\ I_{S_y S_x}(t-1), \dots, I_{S_y S_x}(t-nDt+1), I_{sum}(t-1), \dots, I_{sum}(t-nAt+1)]$$

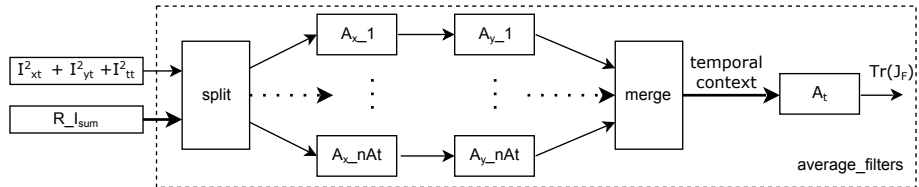
W drugim kroku, w sposób równoległy realizowane są następujące operacje filtr różniczkujący D_y oraz filtr uśredniający S_y . Następnie wykonywane są operacje S_x (dwukrotnie) oraz D_y . W trzecim etapie obliczane są różniczki po czasie D_t i D_{tt} . Niezbędne dane z pamięci RAM opóźniane są w modułach $DEL_{D_y S_x}$, $DEL_{S_y D_x}$ oraz $DEL_{S_y S_x}$. Dodatkowo opóźniane są dane z RAM potrzebne na etapie uśredniania przestrzennego $DEL_{I_{sum}}$. Wyjściem z modułu są wartości I_{xt} , I_{yt} , I_{tt} oraz dane z pamięci RAM oznaczone prefiksem R_{\cdot} .



Rysunek 3: Schemat modułu obliczającego ślad macierzy ($I_{xt}I_{ty}I_{tt}$)

2.3.3. Uśrednianie

Schemat modułu **average_filters** przedstawiono na rysunku 4. Wejście stanowi suma I_{sum} dla bieżącej ramki oraz $nAt-1$ wartości I_{sum} z poprzednich ramek (prefiks R_{\cdot}). W pierwszym kroku dane zostają rozdzielone na nAt wartości. Następnie dla każdej z nich wykonywane są operacje A_x oraz A_y . W kolejnym etapie następuje ich połączenie w temporal context dla którego wykonywana jest operacja A_t . W wyniku otrzymuje się ślad macierzy $Tr(J_F)$.



Rysunek 4: Moduł uśredniający A_x, A_y, A_t

3. Fuzja metody Flux Tensor i Split Gaussian Models

Schemat wykonywanych operacji został zaprezentowany na rysunku 5. Obraz analizowany jest równoległe za pomocą trzech modeli. Pierwszy jest oczywiście model *Flux Tensor*, drugi to model tła reprezentowany przez zmienną liczbę rozkładów Gaussa, natomiast trzecim jest model pierwszoplanowy składający się z jednego rozkładu (przykładowe maski binarne wszystkich trzech modeli zostały przedstawione na rys. 6).

W pierwszej kolejności pod uwagę brany jest wynik z dwóch pierwszych modułów. Na tym etapie może wystąpić jeden z czterech przypadków. Oznaczmy przez F maskę definiującą przynależność piksela do tła lub pierwszego planu. Jest to maska binarna, gdzie 1 oznacza piksel reprezentujący obiekt pierwszoplanowy, a 0 to piksel reprezentujący tło. Do jej wyznaczenia służy następujący wzór:

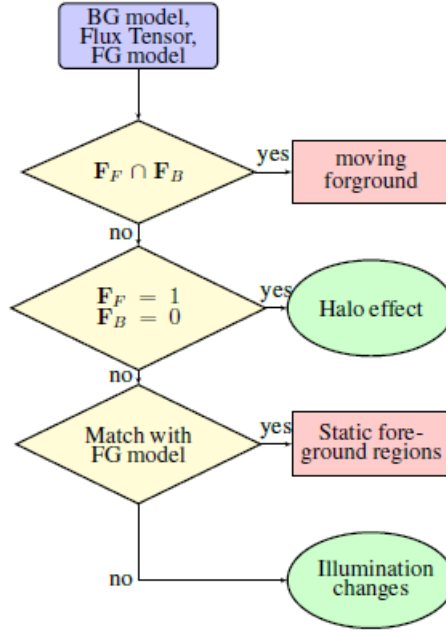
$$F = \begin{cases} 1, & \text{gdy } F_F = 1 \wedge F_B = 1 \\ 0, & \text{gdy } F_B = 0 \\ F_S, & \text{gdy } F_F = 0 \wedge F_B = 1 \end{cases} \quad (6)$$

gdzie:

F_F – maska otrzymana z modelu *Flux Tensor*

F_B – maska otrzymana z modelu tła

F_S – maska definiująca obszar statyczny



Rysunek 5: Schemat operacji algorytmu *FTSG*



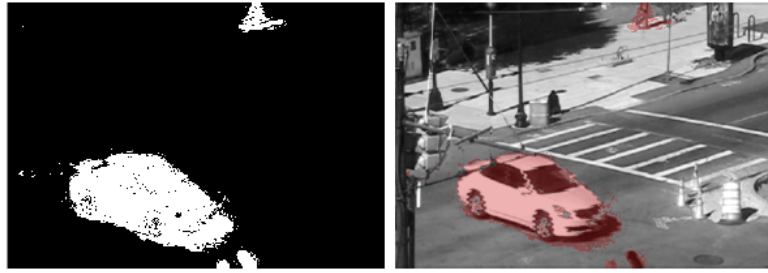
Rysunek 6: Przykładowe maski binarne otrzymane odpowiednio dla modelu tła, metody *Flux Tensor* i modelu pierwszego planu

Jeżeli piksel zostanie przez obie metody tak samo sklasyfikowany, to jego ostateczna przynależność jest oczywista. W takim przypadku piksel zostaje uznany za obiekt ruchomy ($F_F = 1 \wedge F_B = 1$) bądź statyczne tło ($F_F = 0 \wedge F_B = 0$). Na rysunku 7 została przedstawiona część wspólna masek binarnych obu modułów oraz zidentyfikowany obiekt.

Dodatkowa analiza ma miejsce w przypadku, gdy wyniki obu modeli są sprzeczne. Jeżeli metoda *Flux Tensor* oznaczy próbkę jako pierwszy plan, natomiast model tła jako element tła, to końcowym werdyktem jest zaklasyfikowanie piksela jako tło. Jest to logicznie uzasadnione, ponieważ *Flux Tensor* ma tendencję do zaznaczania obiektów o trochę większym obszarze niż są one w rzeczywistości (tzw. efekt halo). Dodatkowo metoda bazująca na rozkładach Gaussa bez problemu wykrywa obiekty będące w ruchu, zatem takie działanie nie powoduje utraty istotnych danych. Taki przypadek został przedstawiony na rys. 8.

W przypadku odwrotnym, czyli $F_F = 0 \wedge F_B = 1$, wykorzystywany jest model pierwszoplanowy i na jego podstawie wyznaczane jest ostateczne dopasowanie. Jeżeli dla zdefiniowanego w tym modelu rozkładu Gaussa piksel wejściowy spełnia warunek (??) to zostaje zakwalifikowany jako element statyczny pierwszego planu ($F_S = 1$). Model pierwszoplanowy jest aktualizowany tylko wtedy, gdy prawdziwy jest warunek $F_F = 0 \wedge F_B = 1 \wedge F_S = 1$. W przeciwnym wypadku ma miejsce aktualizacja modelu tła. Przykład detekcji obszaru statycznego został przedstawiony na rys. 9.

Taka klasyfikacja pozwala wyeliminować większość szumów powstałych w modelu tła, między innymi na skutek zmian oświetlenia i innych zakłóceń (przykładowe szумы w modelu tła zostały przedstawione na rys. 10).



Rysunek 7: Detekcja ruchomego obiektu



Rysunek 8: Przykład tzw. efektu halo

Należy przy tym pamiętać, że współczynnik k w podanym wzorze jest inny dla obu modeli. Istotne jest także, że w niektórych przypadkach obiekt statyczny może być w rzeczywistości odsłoniętym tłem przez właśnie ruszający obiekt (np. samochód odjeżdżający z parkingu). Problem ten jest częściowo eliminowany przez zastosowanie różnych stałych uczenia się (α) w obu modelach wykorzystujących rozkłady Gaussa.

4. Detekcja obiektów zatrzymanych

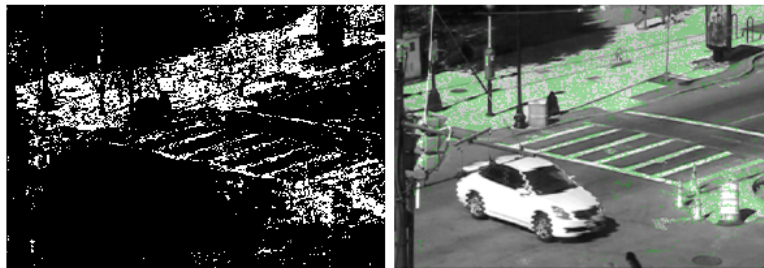
W poprzednim podrozdziale zostało wspomniane, że model pierwszego planu klasyfikuje jako obiekt statyczny także obszar odsłoniętego tła. Jest to wspominany wielokrotnie efekt „duchów”. Może on zostać częściowo wyeliminowany poprzez dobranie odpowiedniego parametru α . Niestety, ze względu na specyfikę metod wykorzystujących rozkłady Gaussa, nie jest możliwa jego całkowita eliminacja. W tym celu został opracowany dodatkowy mechanizm rozróżniania rzeczywistych obiektów zatrzymanych od odsłoniętego tła.

Opisywana metoda wykorzystuje trzy obrazy: wejściowy, model tła oraz maskę binarną otrzymaną z modelu pierwszego planu. Kolejnym krokiem jest wykrycie krawędzi (metodą Sobela) na wszystkich trzech obrazach. Ostatnią czynnością jest porównanie i analiza rezultatów. Jeżeli kształt krawędzi na obrazie wejściowym i wygenerowanym pierwszym planie jest podobny, obiekt może zostać zinterpretowany jako statyczny (tj. zatrzymany obiekt). Natomiast, jeśli krawędzie w modelu tła oraz pierwszego planu są do siebie zbliżone, to wykryty obiekt jest klasyfikowany jako odsłonięte tło. Na rysunkach 11 i 12 zostały przedstawione przykłady ilustrujące oba przypadki.

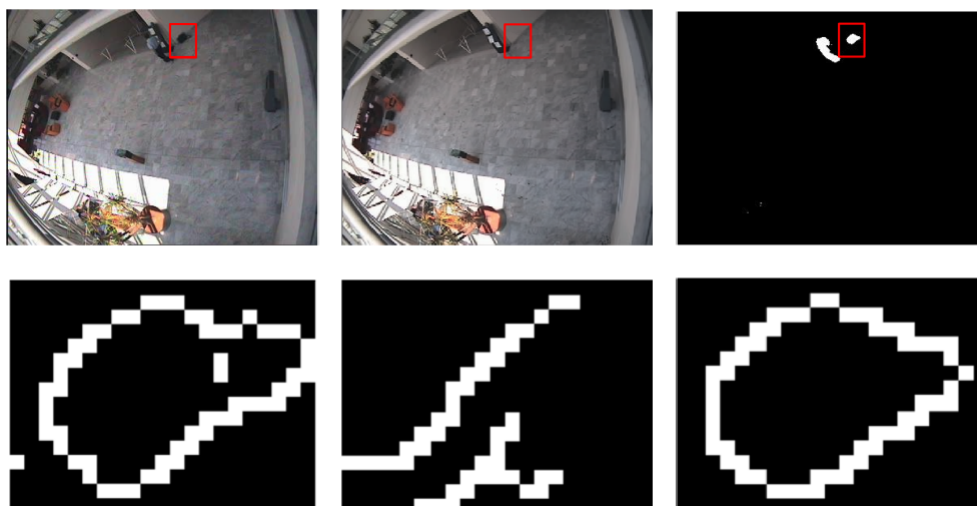
Na rysunku 11 umieszczone zostały w pierwszym rzędzie kolejno: obrazek wejściowy, model tła oraz maska otrzymana przez model pierwszoplanowy. W rzędzie drugim znajdują się obrazy po zastosowaniu detekcji krawędzi. Przedstawiają one zbliżenie na obszar, który został zakwalifikowany jako statyczny przez model pierwszoplanowy i konieczne jest zastosowanie algorytmu detekcji. Jak łatwo zauważyć krawędzie modelu pierwszoplanowego pokrywają się z krawędziami obrazu wejściowego. Zatem obszar ten jest obiektem statycznym. Rysunek 12 przedstawia analogiczną sytuację, jednak w tym przypadku krawędzie modelu tła i modelu pierwszoplanowego są zbliżone, zatem obszar ten należy interpretować jako odsłonięty fragment tła. Takie działanie pozwala wyeliminować obszary błędnie zakwalifikowane, na podstawie modelu pierwszego planu, jako tło.



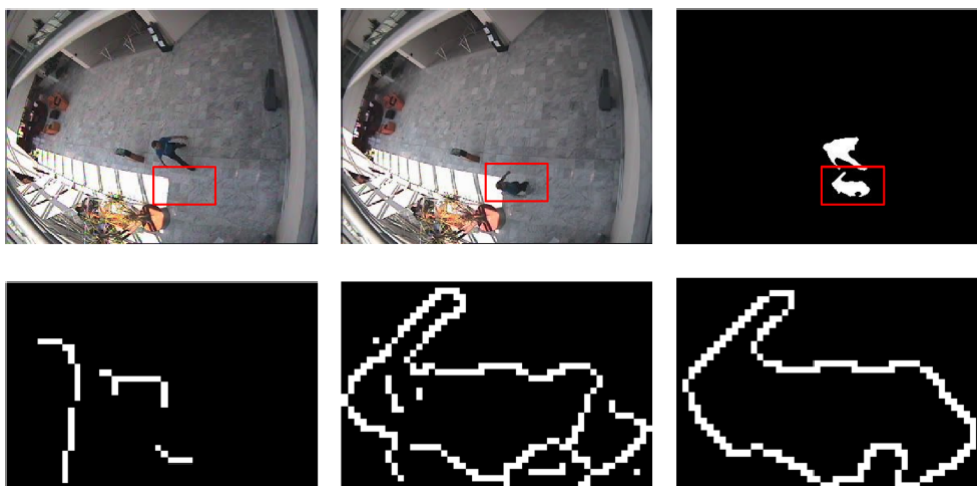
Rysunek 9: Detekcja obszaru statycznego



Rysunek 10: Szumy powstałe w modelu tła



Rysunek 11: obiekt zatrzymany



Rysunek 12: odsłonięte tło