

Rozdział 3.9

Zarządzanie złożonością scenerii na bazie reguł logiki rozmytej

Gabriyel Wong, Nanyang Technological University
gabriel@gmail.com

Jianliang Wang, Nanyang Technological University

Jest rzeczą zrozumiałą, że projektując grę komputerową, należy mieć na uwadze ograniczoną moc obliczeniową dostępnych procesorów potencjalnej platformy docelowej i wynikające stąd ograniczenia w zakresie złożoności renderowania grafiki. Autorzy grafiki trójwymiarowej często korzystają z edytorów sprzężonych wprost z silnikami graficznymi, zyskując w ten sposób możliwość sprawdzenia, czy złożoność ta daje się pogodzić (na konkretnym sprzęcie) z żadaną częstotliwością wyświetlania ramek. Podejście to, które stało się już normą w stosunku do gier rozgrywanych przez pojedynczego gracza — głównie dzięki dużej przewidywalności poczynił tego ostatniego — staje się jednak prawdziwym wyzwaniem w odniesieniu do środowisk dynamicznych, w scenerii współdzielonej przez wielu graczy. W niniejszym rozdziale opiszemy wychodzące naprzeciw temu wyzwaniu podejście oparte na logice rozmytej, umożliwiające pogodzenie efektywności wyświetlania z jakością generowanego obrazu.

Koncepcja

W obliczu ustalonej mocy dostępnych zasobów sprzętowych podstawowym celem grafiki czasu rzeczywistego jest maksymalizacja ostrości (jakości) obrazu odtwarzanego z określoną częstotliwością ramek (wydajność). Mimo iż powszechnie znane rozwiązania w tym zakresie, w rodzaju kontroli poziomu szczegółowości (ang. *Level of Details* — *LOD*) czy redukcji zbioru widocznych obiektów, przyczyniają się wydatnie do redukcji geometrycznej złożoności obrazu, to niezaprzeczalnie stopień tej redukcji zależny jest od konkretnej postaci scenerii i rozwiązania te nie najlepiej nadają się do ogólnej kontroli obciążenia silnika renderującego.

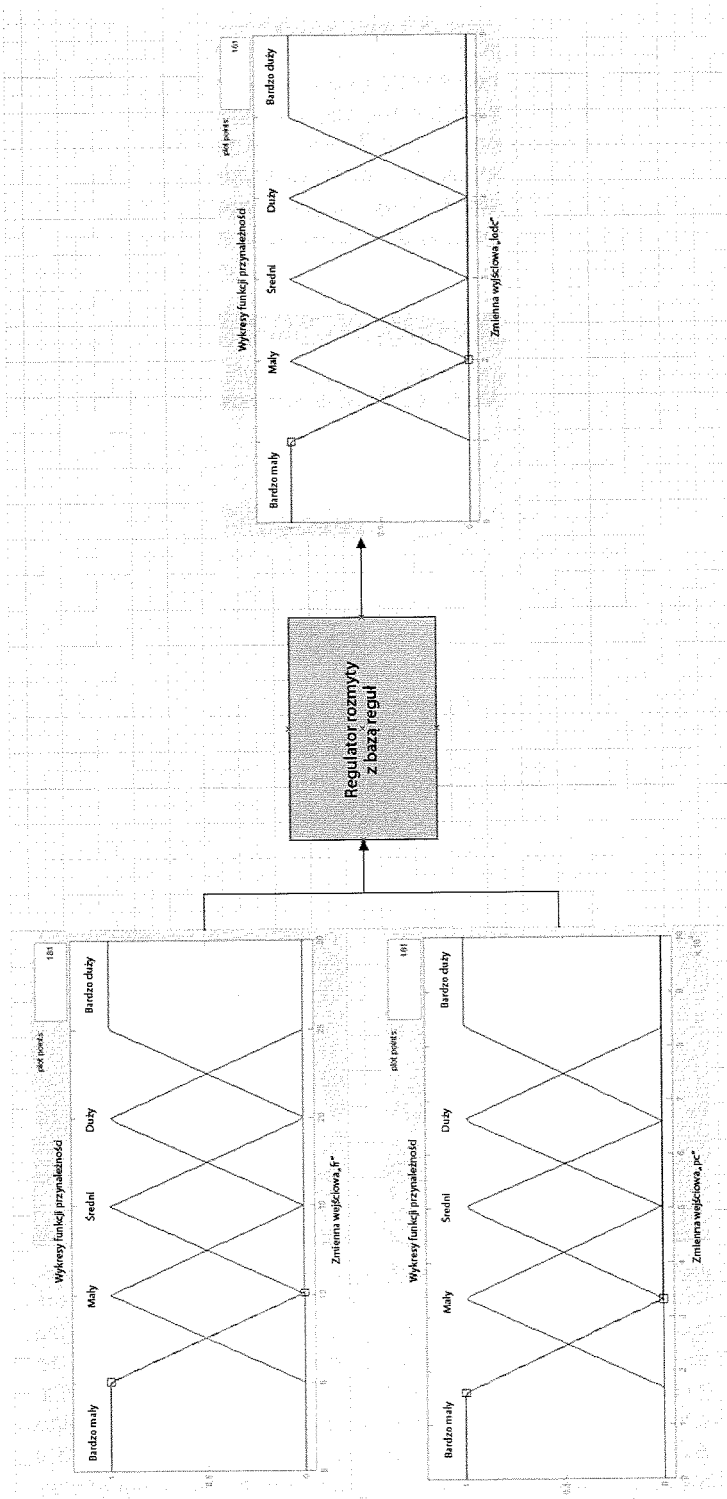
W niniejszym rozdziale nie proponujemy niczego w rodzaju nowego algorytmu unikalnej heurystyki redukcji geometrycznej, prezentujemy w zamian nowatorskie podejście do zarządzania poziomem szczegółowości sceny (LOD) bazujące na zasadach logiki rozmytej (ang. *fuzzy logic*). Utożsamiając proces renderowania grafiki z procesem sterującym, można ustanowić coś na kształt globalnego kontrolera (regulatora) dopływ przepływu mocy „produkcyjnej” (obliczeniowej) do poszczególnych elementów sceny. W połączeniu z mechanizmem sprzężenia zwrotnego kontroler ten, wyposażony w wbudowane reguły wywodzące się z wiedzy projektanta, zdolny jest regulować działanie procesorów drogą rozstrzygnięcia wspomnianego kompromisu „jakość-wydajność”.

Regulacja rozmyta

Zaadaptowanie reguł logiki rozmytej [FIDE1] do potrzeb kontrolowania renderowania daje ze sobą rozliczne korzyści. Przede wszystkim, w przeciwieństwie do konwencyjnych metod sterowania, logika rozmyta dostarcza metodologii szybkiej i wygodnej. Sam proces renderowania jest wystarczająco złożony, by jego ściśle modelowanie uważać za niezwykle wysoce niebanalne; techniki rozmyte dostarczają alternatywy dla modelowania matematycznego w postaci reguł wejścia-wyjścia budowanych na bazie wiedzy projektanta. Implementacja tych reguł upraszcza się do postaci instrukcji warunkowych IF-AND-THEN organizujących relacje między danymi wejściowymi a wyjściowymi. „Lingwistyczne” reguły czyni je bardziej przystępnymi zarówno dla projektantów, jak i użytkowników, bowiem wywodzą się one z opisów formułowanych w języku naturalnym, a nie z zawiłych reguł matematycznych. Ponadto logika rozmyta daje możliwość bardziej efektywnego opisywania wszelkiego rodzaju zjawisk nieliniarnych zachodzących w systemie, które ogólnie pojmowany jest przez projektanta w sposób bardziej intuicyjny. Wszelkie to sprawia, że logika rozmyta staje się atrakcyjną opcją w zakresie podstaw sterowania skomplikowanymi procesami, z renderowaniem grafiki w czasie rzeczywistym na czele.

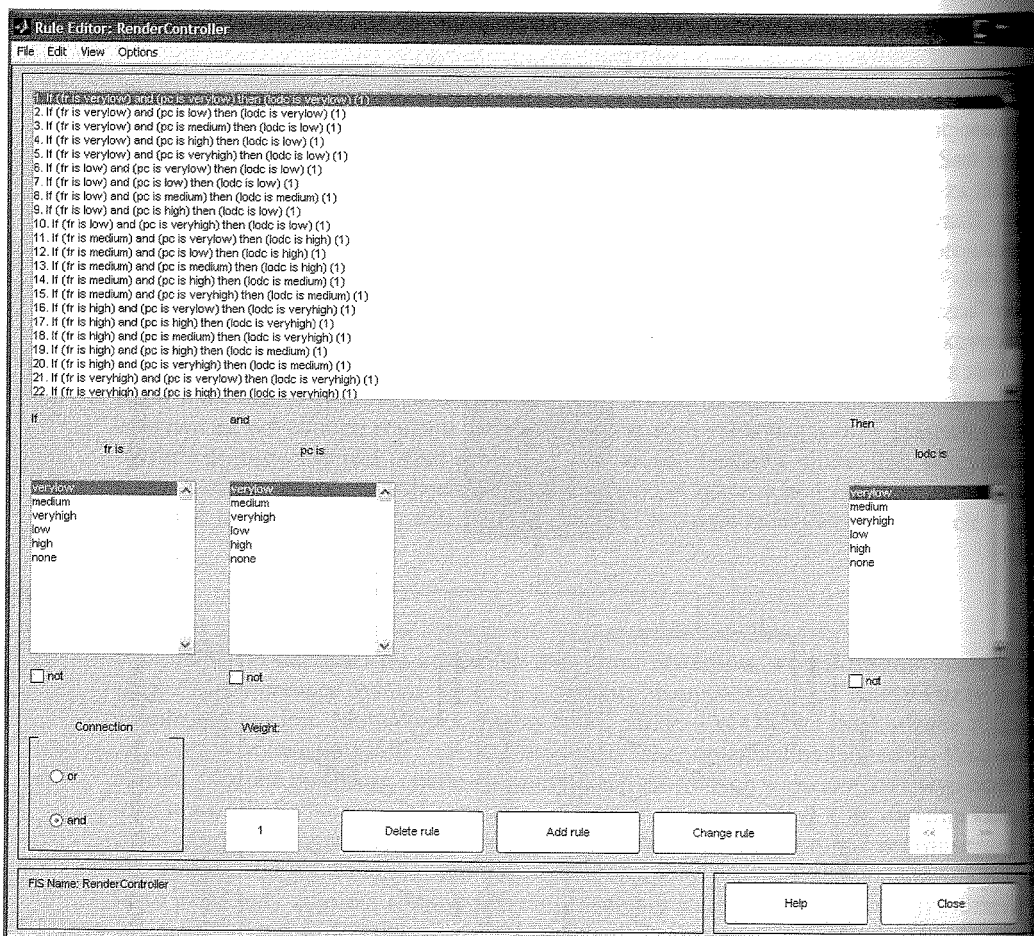
Narzędzia

Istnieje wiele narzędzi do konstruowania tzw. systemów wnioskowania rozmytego (ang. *Fuzzy Inference Systems* — FIS), z bardziej znanych wymienić należy Fuzzy Logic Toolbox z pakietu MathWorks, FuzzyTech[®] firmy INFORM i pochodzące z firmy Apronix środowisko FIDE[™] (Fuzzy Inference Development Environment). Wszystkie oferują podobne podejście, oparte na identyfikacji zmiennych wejściowych i wyjściowych oraz konstrukcji zbioru reguł odzwierciedlających relacje między tymi dwiema kategoriami. Zmienne wejściowe i wyjściowe FIS są obiektami „rozmytymi” w tym sensie, że odzwierciedlają wynik mapowania prostych pojęć języka naturalnego (stanowiących de facto przykłady zbiorów rozmytych) na przedziały wartości liczbowych; wspomniany „mapowanie” ma postać funkcji, zwanej funkcją przynależności (ang. *membership function*), określającej prawdopodobieństwo przynależności konkretnej wartości do konkretnego zbioru rozmytego (czyli do konkretnej kategorii pojęć). Na rysunku 3.9.1 widzimy przykłady mapowania różnych poziomów szczegółowości renderowania (BARDZO MAŁY, MAŁY, ŚREDNI, DUŻY, BARDZO DUŻY) na konkretną wartość, jaką może być liczba pikseli wyświetlanych na obszarze 1 mm² powierzchni wielokąta, czy też liczba wielokątów uwzględnianych w renderowaniu. Wymienione powyżej narzędzia oferują różne postaci funkcji przynależności — alternatywą dla funkcji „piłokształtnej”, widocznej na rysunku 3.9.1, może być np. krzywa Gaussa.



Rysunek 3.9.1. Funkcje przynależności dla zmiennych wejściowych i wyjściowych systemu FIS

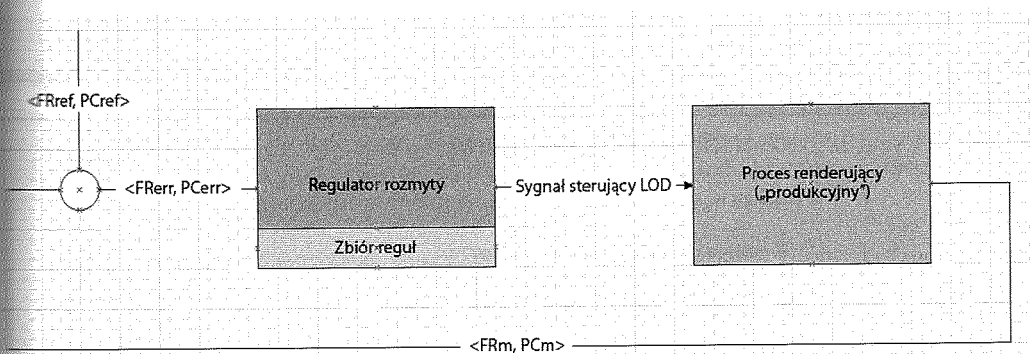
Zadaniem zbioru (bazy) reguł jest opis (w oparciu o wiedzę projektanta) między „poprzednikami” (jakimi są warunki IF) a „konsekwencjami”, czyli en Wysokiej jakości narzędzia do projektowania FIS powinny oferować edytory konstruowanie i edytowanie (dostosowywanie) takich reguł; zrzut ekranu jednego z edytorów widoczny jest na rysunku 3.9.2. Należy przy tym zaznaczyć, że skomponowanie nie najbardziej nawet wyszukanego zbioru reguł nie jest bynajmniej tożsame z konstruowaniem systemu FIS: reguły te muszą zostać poddane serii testów (i przypuszczalnie wielokrotnym modyfikacjom), zanim ostatecznie zweryfikowana zostanie ich użyteczność.



Rysunek 3.9.2. Przykładowy zbiór reguł wiążących zmienne wejściowe systemu FIS ze zmiennymi wyjściowymi

Projekt systemu

Na rysunku 3.9.3 widoczny jest ogólny schemat systemu wykorzystującego regulator rozmyty (ang. *Fuzzy Logic Controller* — *FLC*). Regulator ten sprzężony jest z procesem renderowania („produkcyjnym”), dzięki czemu uzyskuje informacje o liczbie widocznych na ekranie wielokątów i częstotliwości wyświetlania ramek. Niezwykle istotną kwestią związaną z projektowaniem takiego systemu jest jego konfigurowalność: preferowana częstotliwość odtwarzania ramek i maksymalna liczba renderowanych wielokątów to jedne z najważniejszych parametrów regulujących funkcjonowanie aplikacji. Wartości tych parametrów konfrontowane są z informacją otrzymywaną w wyniku wspomnianego sprzężenia, a wynik tej konfrontacji przekazywany jest w postaci zmiennych wejściowych do regulatora.



Rysunek 3.9.3. System renderujący wykorzystujący regulator rozmyty

Regulator rozmyty można potraktować jak czarną skrzynkę, obliczającą wartości zmiennych wyjściowych na podstawie wartości zmiennych wejściowych. Informacja wyjściowa regulatora jest następnie przekazywana do procesu renderującego jako sygnał sterujący wyborem LOD, w konsekwencji czego może zmienić się obciążenie silnika graficznego w związku z renderowaniem wielokątów.

LOD jest podstawowym mechanizmem, za którego pośrednictwem FLC steruje obciążeniem systemu renderującego; sam FLC jako taki nie reprezentuje wprawdzie żadnego algorytmu ani heurystyki związanej z redukcją złożoności geometrycznej, ale jego działanie ukierunkowane jest właśnie na tę redukcję poprzez sprzężenie z procesem renderującym. W książce [Luebke03] znajdują czytelnicy obszerny przegląd różnych technik LOD, zagadnień pokrewnych i ich zastosowań w aplikacjach graficznych czasu rzeczywistego. W tym miejscu ograniczymy się tylko do stwierdzenia, że do najpopularniejszych technik LOD należą techniki dyskretna i ciągła. Istotą techniki dyskretniej jest różnicowanie w czasie liczby uwzględnianych wielokątów, przy niezmiennym (predefiniowanym) sposobie renderowania konkretnego wielokąta; kryterium ustalania wspomnianej liczby może być odległość („głębokość”) poszczególnych wielokątów lub ich zagęszczenie na ekranie. W efekcie otrzymujemy metodę nieskomplikowaną i łatwą w implementacji, jednakże jej skuteczność uwarunkowana jest w dużym stopniu zdolnościami projektanta-grafika, bowiem zbyt duże różnice między kolejnymi poziomami mogą prowadzić do wyraźnych zakłóceń w płynności odtwarzania. W metodzie

ciągłej unika się tego problemu za pomocą „wygładzania” procesu przełączania mów poprzez *stopniową* redukcję szczegółów w miarę wzrostu odległości obiektu. Dodatkowo zastosowanie „siatkowego” renderowania o zróżnicowanych rozdzielczościach także przyczynia się do zmniejszenia złożoności geometrycznej, bez zauważalnych negatywnych konsekwencji, nawet przy małej liczbie wielokątów.

System opisywany w niniejszym rozdziale wykorzystuje ciągły mechanizm LOD dla trójwymiarowych obiektów świata wirtualnego. Wyjście regulatora — zmienna sterująca — pełni rolę mnożnika dla ustalenia LOD i mnożnik ten jest czynnikiem dodatkowy, oprócz wspomnianej już odległości wielokąta od obserwatora. Dzięki owemu dodatkowemu czynnikowi można uzyskać inną niż liniową — kładniczą — zależność liczby wielokątów redukowanych (ignorowanych) od odległości. Fragment kodu na listingu 3.9.1 pokazuje, jak można coś takiego implementować.

Listing 3.9.1. Nieliniowy zanik poziomu szczegółowości uzyskany dzięki sprzężeniu silnika renderującego z regulatorem rozmytym

```
void setCLODFalloff(float falloffMultiplier)
{
    // Dla każdego węzła siatki w scenarii
    {
        // Obliczenie odległości obiektu od kamery
        // Wyrażenie odległości w kategoriach względnych
        float fRatio = (fDepth-fNear)/(fFar-fNear);

        // Gdy żądamy nieliniowego zaniku szczegółowości
        fRatio = Mathf::Pow(fRatio, 2.0f);

        int newIndexCount = (int)(maxIndex * fRatio * (1.25/falloffMultiplier));

        CLODObject->SetIndexCount() = newIndexCount;

        // ...
    }
}
```

Wykorzystując ciągłą technikę redukcji LOD, można osiągnąć efekt „gładkiego” przejścia między kolejnymi poziomami szczegółowości, bez niepożądanych artefaktów widocznych w czasie wyświetlania animacji. A to spełnia warunki podstawowej zasady projektowania systemów renderujących, zgodnie z którą najważniejsze jest zachowanie poziomu jakości obrazu. Opisywany system skonstruowany został za pomocą narzędzi *open-source*; regulator rozmyty zaimplementowano na bazie biblioteki FIS (w języku C++ dostarczonej przez Matlab®).

Zastosowanie do gier

Fuzzy Logic Toolbox z pakietu MathWorks [Matlab1] stanowić może doskonały punkt wyjścia jako narzędzie do konstrukcji systemu FIS wykorzystywanego w grach komputerowych. Przy tworzeniu tego systemu należy jedynie pamiętać o dwóch następujących rzeczach.

- Tworzony system wyposażony jest w predefiniowane reguły wiążące ze sobą obciążenia w zakresie renderowania wielokątów i częstotliwość odtwarzania ramek. Reguły te wydają się fundamentalne, co jednak nie wyklucza ich optymalizacji pod kątem odmiennych zastosowań i różnych platform sprzętowych.
- Konieczne jest uzyskanie pewnego punktu odniesienia, w postaci rezultatów początkowych testów systemu, na potrzeby przyszłego dostosowywania zbioru reguł, aż do uzyskania zadowalających wyników.

Pakiet MathWorks zawiera wygodne narzędzie do prototypowania aplikacji opartych na logice rozmytej — jest nim środowisko Simulink® [Matlab2]. Kreowanie aplikacji w tym środowisku sprowadza się podczas tworzenia odrębnych bloków konstrukcyjnych (ang. *building blocks*) i łączenia ich następnie w spójny system. Nawet prosty plik modelu Simulink ilustruje koncepcję renderowania kontrolowanego przez FLC. Możliwe jest używanie różnych generatorów sygnałów symulujących zróżnicowane obciążenie procesora — zróżnicowane wskutek różnych czynników zewnętrznych, takich jak ruch kamery, kreowanie nowych encji czy włączanie do scenerii nowej kategorii obiektów. Alternatywą dla użycia generatorów sygnałów wejściowych jest zastąpienie abstrakcyjnego procesu „produkcyjnego” rzeczywistą aplikacją (grą) w celu oceny prawdziwej wydajności FLC; umożliwiają to tzw. S-funkcje służące do kreowania „otoczek” (ang. *wrappers*) dla aplikacji zewnętrznych.

Warto również zauważyć, że nie sposób utworzyć FLC, który zdolny będzie do uporańia się z obciążeniem przekraczającym możliwości platformy sprzętowej. Co prawda, projektanci i użytkownicy mają możliwość określania preferowanej częstotliwości odtwarzania ramek i maksymalnej liczby wielokątów (co pokazano na rysunku 3.9.3), lecz wartości te nie mogą wykraczać poza zagregowaną moc obliczeniową wspomnianej platformy. Stąd zalecenie, by na początku przeprowadzić testy umożliwiające ocenę tej mocy i nadanie rozsądnych wartości wspomnianym parametrom.

Obciążenia

W trakcie dotychczasowych rozważań przez cały czas przyjmowaliśmy milcząco założenie, iż zapotrzebowanie aplikacji na moc obliczeniową wyznaczone jest głównie przez geometryczną złożoność scenerii. Założenie to jest jednak o tyle uproszczone, iż na typową aplikację składa się znacznie więcej procesów — łączność sieciowa, sztuczna inteligencja, generowanie dźwięku, symulowanie zjawisk fizycznych itp. — także mających swój udział w budżecie czasowym każdej odtwarzanej ramki. Co więcej, nowoczesne rozwiązania sprzętowe (takie jak układy cieniujące — „shadery”) sprawiają, iż związek między liczbą wielokątów a poziomem szczegółowości ich renderowania staje się bardziej skomplikowany, niż wynikałoby to z prostego mapowania pikseli i cieniowania powierzchni piksel po pikselu.

Nie zmienia to jednak ogólnej zasady, że zasady logiki rozmytej mogą być użyteczne w kontekście sterowania różnymi procesami, także niewiele mającymi wspólnego z grafiką i jej renderowaniem. Wystarczy jedynie dostarczyć do FLC informację (w postaci zmiennych wejściowych) o tychże procesach i utworzyć stosowny zbiór reguł wiążący wartości zmiennych wyjściowych z określonymi kombinacjami zmiennych wejściowych, jak to wcześniej opisywaliśmy. Niezaprzeczalnie wymaga to adekwatnej wiedzy projektanta i należytej staranności w zakresie formułowania rzeczonych reguł.

Uwagi implementacyjne

Mimo iż mechanizm ciągłego dostosowywania LOD dostarcza środki dla bardziej precyzyjnego sterowania obciążeniem geometrycznym silnika, to jednak nie należy zapominać, iż same obliczenia związane z wyznaczaniem i przełączaniem LOD także wymagają pewnej obliczeniowej i jako takie mogą mieć negatywny wpływ na płynność odtwarzania sceny. Zalecane jest więc wykonywanie wspomnianych obliczeń nie dla każdej ramki, ale dla dłuższych okresów czasowych. Utrzymywanie niezmiennego LOD w dłuższym czasie ma także tę zaletę, iż duże (zazwyczaj) podobieństwo sąsiednich ramek (wynikające z ich logicznej spójności w kontekście scenarii) można wykorzystywać do optymalizacji ich wyświetlania. Beztrudne przełączanie LOD co kilka ramek (lub nawet w każdej ramce) mogłoby natomiast powodować szkodliwe migotanie obrazu zwłaszcza wtedy, kiedy obiekty poruszałyby się w niewielkiej odległości od „oka” kamery.

Ujemną stroną ciągłej metody dostosowywania LOD jest konieczność zbudowania struktury danych na temat spójności (wierzchołkowej) grafu obiektów 3D przed wysłaniem do silnika polecenia renderowania tychże obiektów. Mimo iż informację tę można uzyskać struując offline, to i tak wpływa negatywnie na produktywność projektantów, zwłaszcza w warunkach częstych modyfikacji i restartów gry. Ponadto, ponieważ same obliczenia związane z ciągłą metodą LOD stanowią znaczące obciążenie dla procesora, wskazane jest stosowanie tej metody jedynie dla wybranych obiektów scenarii.

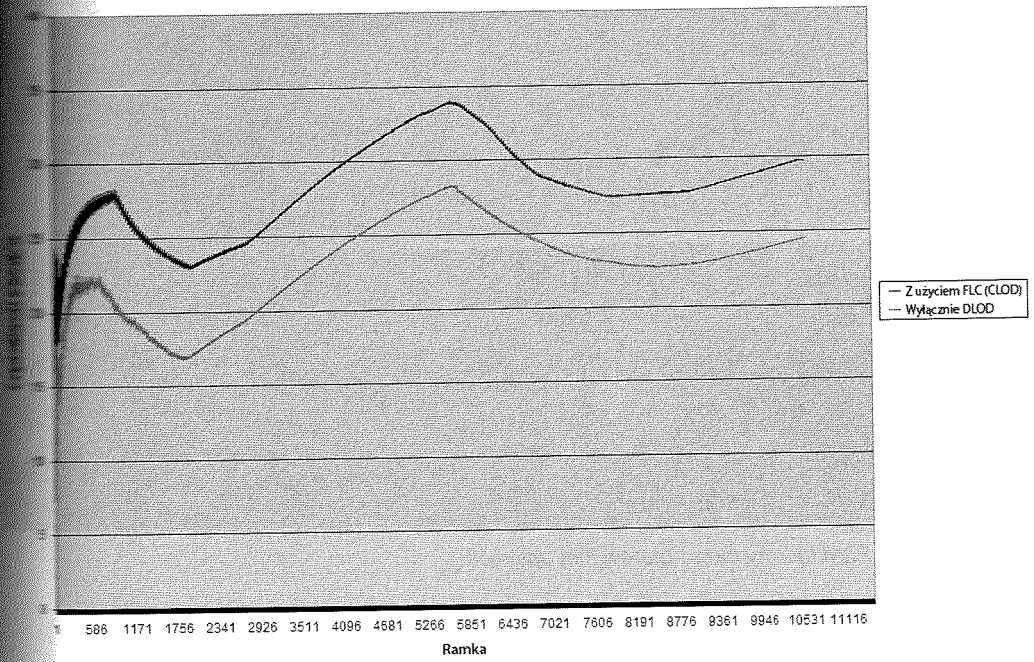
Testy i ich rezultat

Opisywane koncepcje zilustrujemy prostą aplikacją: w przykładowym świecie 3D stworziliśmy na otwartej przestrzeni populację bydła (krów), z których ok. 40% zarządzane jest przez ciągłą metodę dostosowywania LOD (CLOD) — wskaźnik ów określa jednocześnie udział regulatora FLC w zarządzaniu całą scenarią i może być — oczywiście — zmieniony stosownie do konkretnej aplikacji i kontekstu scenarii. Przykładowo złożone obiekty najeżone szczegółami lepiej poddają się kontroli LOD niż licznych, małych i prostych obiektów — te ostatnie wykazują znacznie większe zapotrzebowanie na moc obliczeniową.

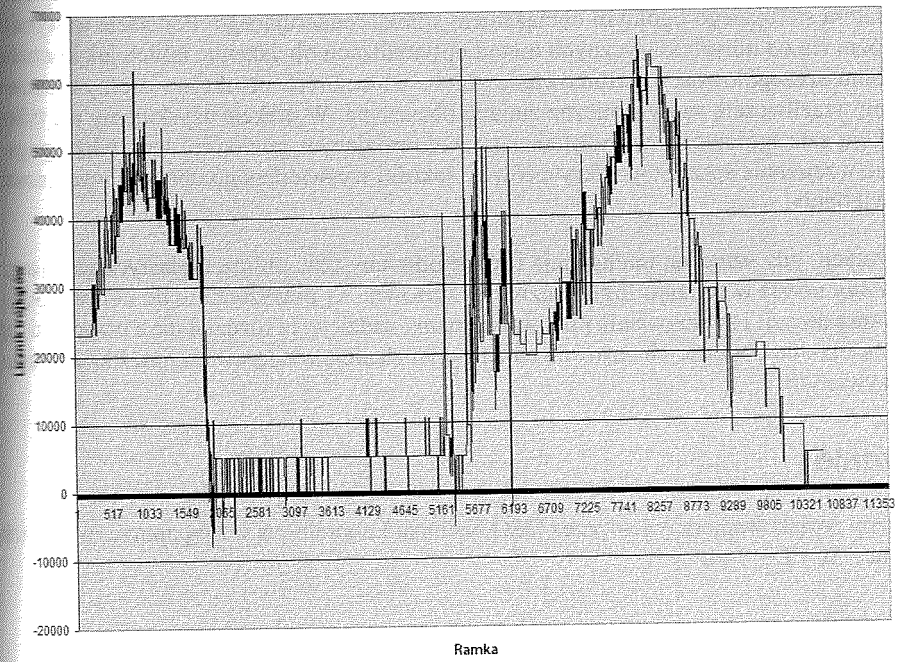
Celem przeprowadzonych testów było zarówno potwierdzenie tezy, że zastosowanie FLC istotnie przyczynia się do poprawy wydajności aplikacji, jak i przetestowanie funkcjonalności oraz efektywności FLC w zakresie sterowania scenarią, zwłaszcza w porównaniu z przypadkiem wykorzystywania jedynie dyskretnej metody LOD (DLOD). Wyniki porównawcze dla obu opisanych podejść widoczne są na rysunku 3.9.4.

Spoglądając na rysunek, można zauważyć, że zastosowanie FLC skutkuje wyższą (o ok. 25% częstotliwością odtwarzania ramek, w porównaniu z systemem wykorzystującym wyłącznie metodę DLOD. Oba systemy poddane zostały na początku testu identycznemu, silnemu obciążeniu wynikającemu z konieczności odtwarzania ramek w tempie 30 na sekundę. Diagram w dolnej części rysunku przedstawia względną różnicę liczby trójkątów w obu systemach dla każdej ramki — ewidentnie system wykorzystujący FLC wykazywał większą zdolność zachowywania szczegółów, mimo iż rzeczywista częstotliwość odtwarzania spadła miejscami do ok. 22 ramek/s. Przewaga ta konsekwentnie utrzymywała się w dalszym ciągu testu, w stopniu zależnym od geometrycznej komplikacji scenarii. Przykładowy zrzut ekranu z testowanej aplikacji widoczny jest na rysunku 3.9.5 oraz na kolorowej ilustracji nr 6.

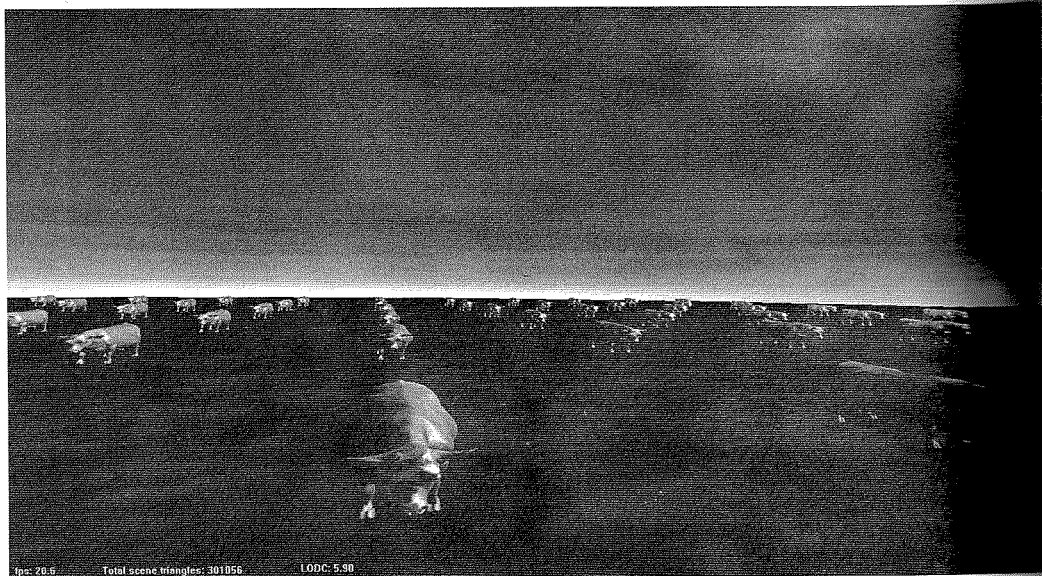
Wykres częstotliwości odtwarzania ramek



Różnica liczby trójkątów



Rysunek 3.9.4. Porównanie efektywności dwóch różnych systemów renderujących



Rysunek 3.9.5. Przykładowy ekran testowanej aplikacji

Konkluzja

W niniejszym rozdziale opisaliśmy zalety regulacji rozmytej zastosowanej do zarządzania złożonością scenarii gier; dzięki tym zaletom możliwe okazało się skonstruowanie regulatora sterującego ogólnym obciążeniem silnika renderującego. Podejście to nie różni się od innych klasycznych metod optymalizacji obciążenia, takich jak kontrola poziomu szczegółowości i redukowanie zbioru widocznych obiektów, istotnym szczegółem: wrażliwością na dynamiczne zmiany obciążenia. W efekcie uzyskujemy nie tylko redukcję szczegółowości wyświetlania na zadowalającym poziomie, lecz także pogodzenie wzajemnie sprzecznych kryteriów — jakości wyświetlanego obrazu i wydajności wyświetlania.

W sytuacji, gdy sztuczna inteligencja i różne odmiany tzw. obliczeń elastycznych (*soft-computing*) w rodzaju logiki rozmytej stosowane są powszechnie do celów podejmowania decyzji i symulowania pola walki, mamy nadzieję, iż rozdział nasz przyczyni się do poszerzenia obszaru ich zastosowań na optymalizację renderowania grafiki w rzeczywistym.

Podziękowanie

Opisany projekt wspierany jest przez DSO National Laboratories w Singapurze, w ramach grantu DSOCL01144.