

**Fehér Krisztián**

# **Hogyan írnj extrém gyors programot?**

**- Bevezetés a CUDA programozásba -**

BBS-INFO Kiadó, 2019.

Minden jog fenntartva! A könyv vagy annak oldalainak másolása, sokszorosítása csak a kiadó írásbeli hozzájárulásával történhet.

A könyv nagyobb mennyiségben megrendelhető a kiadónál:  
BBS-INFO Kiadó, [www.bbs.hu](http://www.bbs.hu) Tel.: 407-17-07

A könyv megírásakor a szerző és a kiadó a lehető legnagyobb gondossággal járt el. Ennek ellenére, mint minden könyvben, ebben is előfordulhatnak hibák. Az ezen hibákból eredő esetleges károkért sem a szerző, sem a kiadó semmiféle felelősséggel nem tartozik, de a kiadó szívesen fogadja, ha ezen hibákra felhívják figyelmét.

ISBN 978-615-5477-78-2  
E-book ISBN 978-615-5477-79-9

Kiadja a BBS-INFO Kft. Budapest  
Felelős kiadó: a BBS-INFO Kft. ügyvezetője  
Nyomdai munkák: Biró Family Nyomda  
Felelős vezető: Biró Krisztián

# TARTALOMJEGYZÉK

<b>1. Előszó</b> .....	7
1.1. A téma aktualitása.....	7
1.2. A könyv felépítése.....	8
1.3. Kinek szól a könyv? .....	8
1.4. A szerzőről .....	9
<b>2. A CUDA madártávlattól</b> .....	<b>10</b>
2.1. Mi az a CUDA?.....	10
2.2. CUDA generációk .....	10
2.3. Mire használják a CUDA-t? .....	11
2.3.1. Mesterséges intelligencia .....	11
2.3.2. Adatközpontok .....	11
2.3.3. GPU felhő.....	11
2.3.4. Dizájn és vizualizáció, fizikai szimulációk.....	11
2.3.5. Autonóm eszközök, önvezető autók.....	12
2.3.6. Játékipar .....	12
2.3.7. Bioinformatika, vegyipar .....	12
2.3.8. Adatkutatás.....	12
2.3.9. Hadiipar .....	12
2.3.10. Automatizált elektronikai tervezés.....	13
2.3.11. Pénzügyi számítások .....	13
2.3.12. Orvosi képalkotás.....	13
2.4. Hogyan programozhatunk CUDA-ban? .....	13
2.4.1. Központi oldal.....	14
2.4.2. Támogatott eszközök .....	14
2.4.3. Fejlesztőkészlet, dokumentáció .....	14
2.5. A gyors adatfeldolgozás alapfeltételei .....	15
2.6. Mennyivel gyorsabbak a CUDA magok? .....	15
2.7. Összegzés .....	15
<b>3. CUDA alapú kártyák bemutatása</b> .....	<b>16</b>
3.1. GeForce, Quadro, Tesla kártyák.....	16
3.2. GeForce.....	16
3.3. Quadro.....	17

3.4.	Tesla .....	17
3.5.	Összegzés .....	17
<b>4.</b>	<b>Fejlesztőrendszer előkészítése .....</b>	<b>18</b>
4.1.	A Visual Studio telepítése .....	18
4.2.	A CUDA Toolkit telepítése .....	18
4.3.	CUDA verziószámozás értelmezése.....	19
4.4.	Meghajtóprogramok .....	19
4.5.	Összegzés .....	19
<b>5.</b>	<b>CUDA lekérdezések.....</b>	<b>20</b>
5.1.	Alapfogalmak .....	20
5.2.	Végrehajtási konfiguráció .....	20
5.3.	Új CUDA projekt létrehozása Visual Studioban .....	21
5.4.	A legegyszerűbb program .....	22
5.5.	CUDA eszközök jelenléte.....	23
5.6.	Eszközadatok lekérdezése .....	23
5.7.	További hasznos segédprogramok .....	26
5.7.1.	bandwidthTest - Memória adatátvitel sebessége .....	27
5.7.2.	topologyQuery - Topológia lekérdezése .....	27
5.8.	Összegzés .....	28
<b>6.</b>	<b>CUDA párhuzamos programozási alapok.....</b>	<b>29</b>
6.1.	Memóriaterületek kezelése, elérése .....	29
6.1.1.	Menedzselte memóriaelérés.....	31
6.2.	Szálak, blokkok és társaik .....	32
6.3.	Párhuzamosítási technikák .....	34
6.4.	Egyszerű példa, CPU változat.....	35
6.5.	GPU kód – Aktuális végrehajtási szál meghatározása .....	36
6.6.	A legegyszerűbb program, kernellel .....	38
6.7.	Párhuzamos végrehajtás blokkok és szálak felhasználásával .....	39
6.7.1.	Több szál megadása.....	39
6.7.2.	Több blokk megadása.....	40
6.7.3.	Szinkronizáljunk! .....	41
6.8.	Kétdimenziós feldolgozás konfigurációja .....	42
6.9.	Összegzés .....	45
<b>7.</b>	<b>Nagyteljesítményű grafika, CPU kóddal .....</b>	<b>46</b>
7.1.	A Direct2D megoldás bemutatása .....	47
7.2.	Pufferelés megvalósítása.....	48
7.3.	2D alakzatok kirajzolása.....	55
7.3.1.	Pont kirajzolása .....	55
7.3.2.	Vonal kirajzolása.....	56
7.3.3.	Háromszög kirajzolása.....	58

---

7.4.	Z-pufferelés.....	61
7.4.1.	Z-pufferelt Pont kirajzolása.....	62
7.4.2.	Z-pufferelt Vonal kirajzolása.....	63
7.4.3.	Z-pufferelt Háromszög kirajzolása.....	64
7.5.	Összegzés.....	65
<b>8.</b>	<b>Extrém grafikus teljesítmény – CUDA-val.....</b>	<b>68</b>
8.1.	Hol gyorsítsunk?.....	68
8.2.	Pufferelés megvalósítása.....	68
8.3.	2D alakzatok kirajzolása.....	69
8.3.1.	Pont kirajzolása.....	69
8.3.2.	Vonal kirajzolása.....	69
8.3.3.	Háromszög kirajzolása.....	70
8.4.	Z-pufferelés.....	72
8.4.1.	Z-puffer inicializálása.....	72
8.4.2.	Z-pufferelt pont kirajzolása.....	72
8.4.3.	Z-pufferelt vonal kirajzolása.....	72
8.4.4.	Z-pufferelt háromszög kirajzolása.....	74
8.5.	Összegzés.....	76
<b>9.</b>	<b>Komplex példa: extrém gyors 3D OBJ modell megjelenítő.....</b>	<b>77</b>
9.1.	Példák.....	84
9.2.	Láthatóság vizsgálata.....	88
9.3.	Fény használata.....	89
9.4.	Textúrák alkalmazása.....	92
<b>10.</b>	<b>Záró gondolatok.....</b>	<b>94</b>
<b>11.</b>	<b>Függelék.....</b>	<b>95</b>
11.1.	Több videokártya használata.....	95
11.2.	CUDA programok teljesítménymérése.....	97
11.2.1.	Általános teljesítménymérés.....	97
11.2.2.	Teljesítménymérés CPU-val.....	97
11.2.3.	Teljesítménymérés CUDA-ban.....	98
11.3.	CUDA programok debuggolása.....	99
11.3.1.	A „csináld magad” módszer.....	99
11.3.2.	Hibakeresés az NSight beépülő modullal.....	99
11.3.3.	CUDA hibák kiírása.....	100
11.4.	Futtatási profilok elemzése.....	100
11.5.	A könyvben használt CUDA függvények.....	102
<b>12.</b>	<b>Ajánlott irodalom.....</b>	<b>103</b>



# 1. Előszó

Izgalmas technológiai áttörésekről olvashatunk nap mint nap, melyek mesterséges intelligenciáról, önvezető autókról és hasonló technológiákról szólnak. A számos úttörő megoldás jelentős hányadának háttérében az NVIDIA vállalat CUDA platformja áll, mely egyedülálló hardveres és szoftveres megoldásaival forradalmi lehetőségeket kínál a jelen és a jövő kutatóinak és a technológia hétköznapi felhasználóinak is.

Ez a könyv a CUDA technológia programozásának alapjairól szól.

## 1.1. A téma aktualitása

Az emberiség által kitermelt digitális adatmennyiség napjainkra gigászi méreteket ölt. A tárolás mellett most már az adatok elemzése, feldolgozása is egyre inkább előtérbe kerül, melynek igényét a mesterséges intelligenciával és a különböző autonóm rendszerekkel kapcsolatos kutatások csak tovább fokozzák.

Kb. 10 évvel ezelőttig érvényben volt az ún. Moore törvény, mely kimondta, hogy a számítógépek számítási kapacitása évente megduplázódik. Ezt a szabályt még az iskolákban is tanították, annyira szemléletes volt. A CPU-k fejlődése kb. 10 évvel ezelőtt azonban megtorpant és azóta lényegében stagnál. Minimális előrelépések történnek csak, mivel a jelenlegi architektúrák és a fizika törvényei már nem teszik lehetővé a számítási teljesítmény további növelését. Moore törvénye nem érvényes többé!

Több kutatás is folyik alternatív, gyökeresen új felépítésű architektúrák kifejlesztése érdekében, de ezek belátható időn belül nem lesznek elérhetőek, pláne bárki által programozhatóak. Megoldásokra márpedig szükség van már most is. Így terelődik a figyelem egyre inkább a videokártyák párhuzamos számítási képességeinek általános célokra történő kihasználására.

Elérkezett a GPU-k korszaka, melyet angol szóval 'GPU computing', vagy 'High Performance Computing' elnevezésekkel is szoktak illetni. Szintén elterjedt kifejezés a GPGPU is (General-Purpose Com-

puting on Graphics Processing Units). A GPU-k a CPU-kkal szemben eleve párhuzamosított, aszinkron adatfeldolgozásra lettek kitalálva. A jelen GPU-k a párhuzamosított adatfeldolgozás terén messze túlszárnyalják a CPU-kat, akár 10 -, vagy 100 szorosán is, nagyon alacsony fogyasztás és nem utolsó sorban alacsony ár mellett. Ezek olyan fegyvertények, amelyek mellett az informatika iránt érdeklődők nem mehetnek el csukott szemmel.

Természetesen a CPU-k sosem fognak eltűnni, de ahol kiugróan nagy számítási teljesítményre lesz szükség, ott a GPU-k használata egyre erőteljesebb lesz. A nagyteljesítményű adatfeldolgozás köré felépülő technológiák a számítástechnika új, napjainkban zajló forradalmát képviselik, ezért mindenképpen ajánlott a folyamatot megismerni: miről szól ez az egész, milyen elvek szerint működik, hol lehet bekapcsolódni ebbe, stb.

## 1.2. A könyv felépítése

A könyv egyedülálló módon egyszerre ismeretterjesztő kiadvány és programozási segédlet is egyben.

A könyv első felében a CUDA technológia történetéről, jelenlegi felhasználási területeiről szólnak.

A könyv második fele a technológiát a gyakorlatban is felhasználni kívánó, a programozás iránt érdeklődő olvasóknak szól.

A függelék további hasznos témaköröket tárgyal.

Minden, a könyvben bemutatott kód korlátozás nélkül felhasználható, bármilyen célra. A fő forráskódok letölthetőek a kiadó weboldaláról is.

A könyv egyes algoritmusainak kidolgozásához a szerző sok inspirációt merített Dmitry V. Sokolov `tinyrenderer` publikus rasterizációs algoritmusából: <https://github.com/ssloy/tinyrenderer>.

## 1.3. Kinek szól a könyv?

Az elsődleges célközönség a C programozásban jártas olvasó, aki rendelkezik tapasztalattal Windows alkalmazások fejlesztésében.

**A WIN32 API, a Direct2D és a Visual Studio mint fejlesztőeszköz alapszintű ismerete elengedhetetlen.** Ha ezek a témák ismeretlenül csengenek, a szerzőnek az ajánlott irodalomban található könyvei segítségével pótolhatóak ezen alapismeretek.

## 1.4. A szerzőről

A szerző hivatásos szoftvertesztelő, minőségbiztosítási tanácsadó, diplomás német irodalmár, a Magyar Térinformatikai Társaság (HUNAGI) egyéni szakértői tagja, független kutató.

Elsődleges szakterülete a digitális grafika programozása, digitális térképalkalmazások készítése és a párhuzamos programozási technikák felhasználási területei.

Kibervédelmi kérdéseket évek óta tudatosan tanulmányoz, több kiberbiztonsággal kapcsolatos előadást is tartott már Magyarországon, sőt vendégoktatóként még oktatási intézményben is. Szakmai munkáját évről évre növekvő érdeklődés és elismerés kíséri.

Gyerekkorában autodidakta módon tanult meg programozni, az évek során számos programozási nyelvvel megismerkedett. megszerzett tudását előszeretettel használja alternatív, kísérleti alkalmazások készítésére, melyek egy része ingyenesen elérhető, sőt vannak köztük nyílt forráskódúak is. A szerző fejleszt Windows desktop, Android és webes környezetekre is.

Tudását igyekszik minél szélesebb körben megosztani másokkal is. Ennek folyamányaként több könyve is megjelent már a hazai könyvesboltokban az elmúlt években, nem egy közülük sikerlisták élére is került.

## 2. A CUDA madártávlattól

### 2.1. Mi az a CUDA?

A CUDA a **Compute Unified Device Architecture** rövidítése, mely magyarul kb. annyit tesz: **Egységesített Számítási Eszközarchitektúra**.

Az NVIDIA a CUDA-t napjainkban **számítási platformként (CUDA-X)** definiálja, ami jól megragadja a technológia lényegét is.

A CUDA első megjelenése 2006 körülre datálható, a GeForce 8800-as kártya megjelenésére, mely a világon elsőként kínált általános célú számítási feldolgozóegységeket, közismertebb nevükön **CUDA magokat**.

Az azóta eltelt időszakban már több mint 500 millió, CUDA technológiát alkalmazó grafikus kártyát adtak el, ami példátlan siker.

A CUDA segítségével általános feladatokat gyorsíthatunk fel, CPU-alapú megközelítésekhez képest jelentős mértékben.

### 2.2. CUDA generációk

Az NVIDIA a CUDA architektúráit hagyományosan valamilyen híres tudósról nevezi el.

A legfontosabb architektúrák elnevezése, időben:

- Fermi
- Kepler
- Maxwell
- Pascal
- Volta
- Turing.

Egy-egy új architektúrát jellemzően pár éves időközökben hoznak ki. A Volta részben kivételt képez, ugyanis ebből az architektúrából csupán néhány felső kategóriás gyorsítókártya került ki. Az összes többi esetén azonban GeForce és Quadro kártyák is rendre készültek.

## 2.3. Mire használják a CUDA-t?

Tekintsük át a legfontosabb szakterületeket, melyeken a CUDA domináns szerepet tölt be!

### 2.3.1. Mesterséges intelligencia

Az MI kutatások sikere belátható időn belül az adatfeldolgozás gyorsaságán áll vagy bukik. Jelenleg az ún. **mély tanulásos** (angolul: **deep learning**) alapú MI kutatások dominálnak, melyek keretében MI-kezdemenyeket tanítanak különböző viselkedésekre, rengeteg mintaadat feldolgozása útján. Éppen ezért belátható, hogy a CUDA az MI kutatások egyik jelentős hajtómotorja.

### 2.3.2. Adatközpontok

A modern adatközpontok felépítése és üzemeltetése a folyamatosan növekvő igények miatt egyre nagyobb kihívást jelent. A kis helyigény, nagyfokú integráltság és az alacsony energiaszükséglet döntő tényezők szoktak lenni. Ezen a területen az NVIDIA Tesla gyorsítókártyái például kimagaslóan jók.

### 2.3.3. GPU felhő

Óriási az igény a változtatható számítási kapacitások bérleti alapon történő felhasználására. Ilyet számos vállalat kínál, például az Amazon is. Segítségével alacsony költségek árán lehet óriási számítási teljesítményeket időszakosan kibérelni.

A másik terület a streaming alapú online videojátékipar, mely a játékosoknak alacsony költségen kínál hozzáférést nagyteljesítményű videokártyákhoz.

### 2.3.4. Dizájn és vizualizáció, fizikai szimulációk

A különböző tervezési feladatoknál nagyon fontos tényező az idő, több szempontból is. Minél gyorsabban előállítható egy-egy terv végleges változata, annál értékesebb egy technológia.

Felbecsülhetetlen értékű például egy olyan technológia, aminek a segítségével például valós időben lehet látni egy még csak tervezési fázisban levő gépjármű megjelenését, vagy akár viselkedését is.