



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI



Instrumente Structurale
2007-2013



Platformă de e-learning și curriculum e-content pentru învățământul superior tehnic

Arhitectura Sistemelor de Calcul

36. Tehnici de optimizare utilizare in calculul de inalta performanta



Memoriei

- Exista mai multe moduri de a imbunatatii performanta cache-urilor si anume
 - Reducerea penalitatilor unui Cache Miss
 - Reducerea ratei de Cache Miss-uri
 - Reducerea timpului de rezolvare a unui Hit
 - Imbunatatirea memoriei
 - Performanta crescuta
 - Cost redus

Optimizari de Compilator



UNIUNEA EUROPEANĂ



GUVERNUL ROMÂNIEI



Instrumente Structurale
2007-2013



- Sunt mai multe moduri in care un cod poate fi modificat pentru a genera mai putine miss-uri
 - Compilatorul
 - Utilizatorul
- Vom analiza un exemplu simplu – initializarea unui vector bidimensional (matrice)
- Vom presupune ca avem un compilator neperformant si vom optimiza codul in mod direct
 - Un compilator bun trebuie sa poata face acest lucru
 - Cateodata insa, compilatorul nu poate face tot ce este necesar

Exemplu: Initializarea unui vector 2-D



```
int a[100][100];
```

```
for (i=0;i<100;i++) {
    for (j=0;j<100;j++) {
        a[i][j] = 11;
    }
}
```

```
int a[100][100];
```

```
for (j=0;j<100;j++) {
    for (i=0;i<100;i++) {
        a[i][j] = 11;
    }
}
```

- Care varianta este mai buna?
 - i,j?
 - j,i?
- Pentru a raspunde la aceasta intrebare, trebuie sa intelegem modul de asezare in memorie al unui vector 2-D

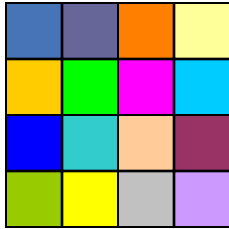
Vectori 2-D in Memorie

- Un vector 2-D static se declara:

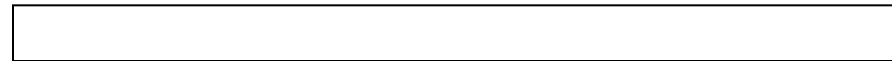
```
<type> <name>[<size>][<size>]  
int array2D[10][10];
```

- Elementele unui vector 2-D sunt salvate in celule **contigue** de memorie
- Problema este ca:
 - Matricele sunt conceptual 2-D
 - Memoria unui sistem de calcul este 1-D
- Memoria 1-D a sistemului este descrisa de un singur numar: adresa de memorie
 - Similar cu numerele de pe axa reala
- Astfel, este necesara o mapare de la 2-D la 1-D
 - Cu alte cuvinte, de la abstractizarea 2-D utilizata in programare, la implementarea fizica in 1-D

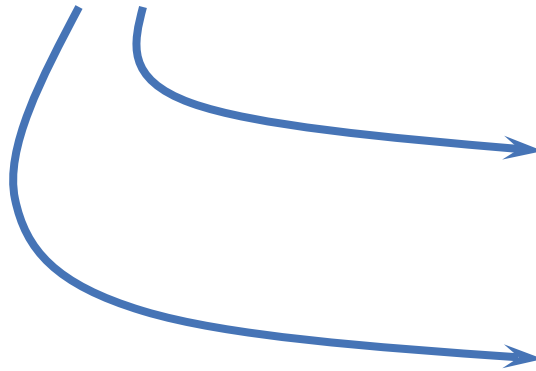
Maparea de la 2-D la 1-D



Matrice $n \times n$ 2-D



Memoria sistemului 1-D



Maparea 2-D la 1-D

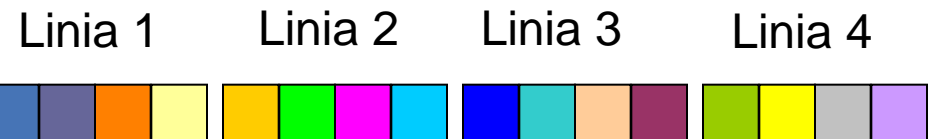
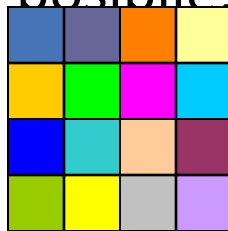


O alta mapare 2-D la 1-D

Exista $n!$ mapari posibile

Row-Major vs. Column-Major

- Din fericire, în orice limbaj sunt implementate maxim 2 din cele $n^2!$ mapări posibile, și anume:



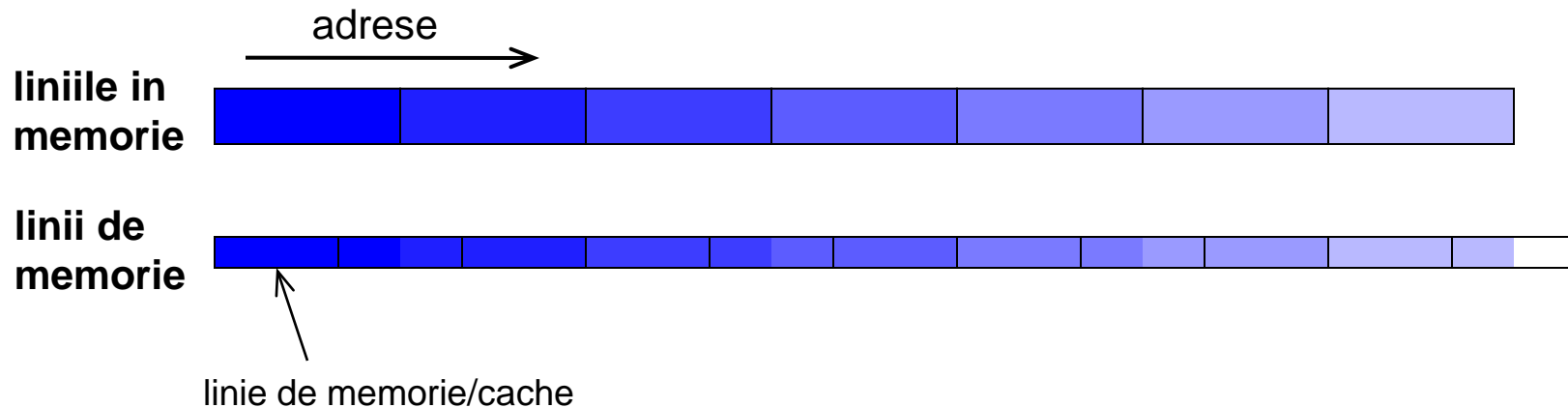
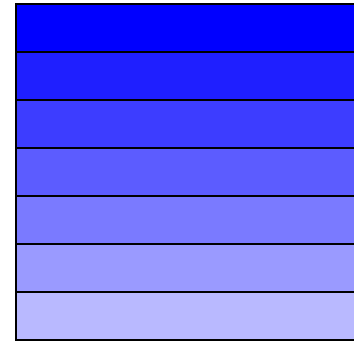
- Row-Major:
 - **Linile sunt stocate continuu**



- Column-Major:
 - **Coloanele sunt stocate continuu**

Row-Major

- Row-Major este utilizat de C/C++

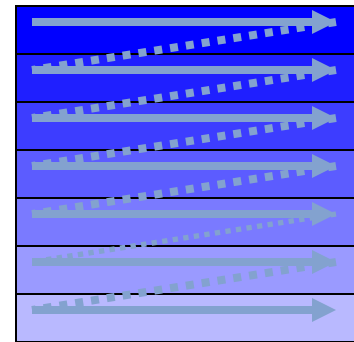


- Elementele matricii sunt stocate contiguu in memorie

Row-Major

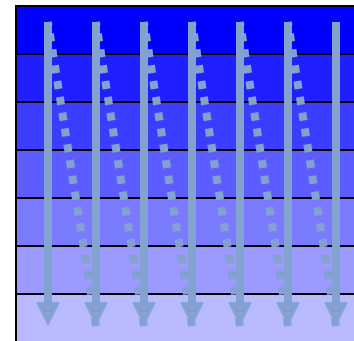
- C/C++ utilizează Row-Major
- Prima implementare (i, j)

```
int a[100][100];  
for (i=0;i<100;i++)  
    for (j=0;j<100;j++)  
        a[i][j] = 11;
```



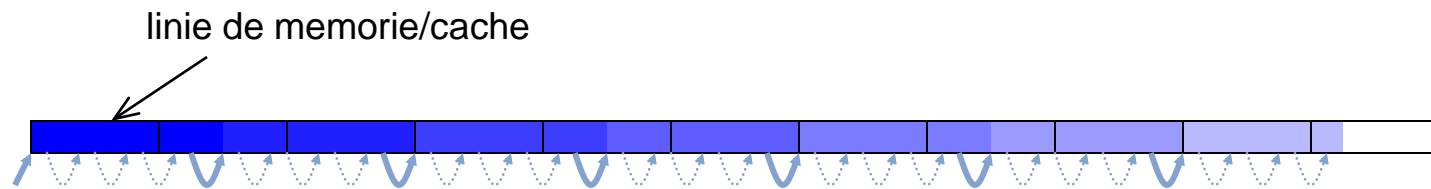
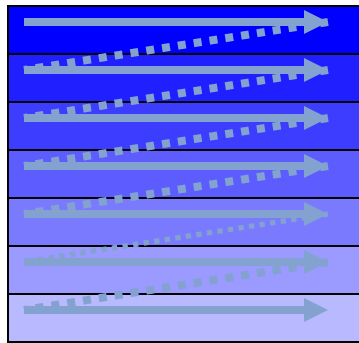
- A doua implementare (j, i)

```
int a[100][100];  
for (j=0;j<100;j++)  
    for (i=0;i<100;i++)  
        a[i][j] = 11;
```

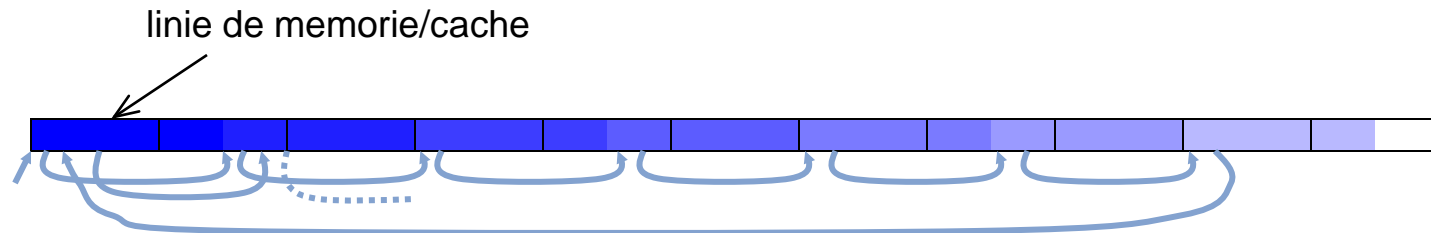
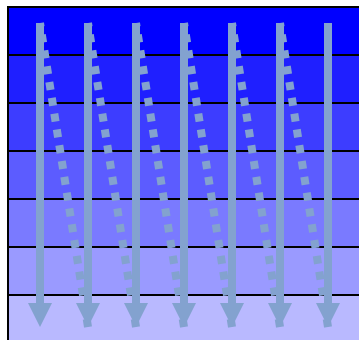


Definiții: Numarul de Miss-uri

- Matricea $n \times n$ 2-D array,
- Fiecare element are e bytes,
- Dimensiunea liniei de cache este b bytes



- Se obține un miss la fiecare linie de cache: $n^2 \times e / b$
- Dacă avem es: n^2 accesuri la memorie (toată matricea)
- Rata de miss-uri este: e/b
- Exemplu: Miss rate = 4 bytes / 64 bytes = **6.25%**
 - În afara cazului în care vectorul este foarte mic

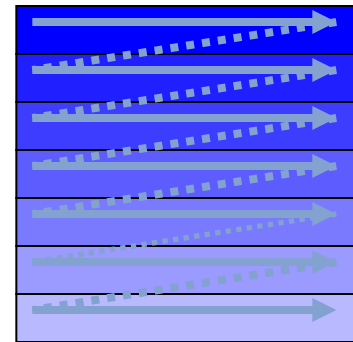


- Avem un miss la fiecare acces
- Exemplu: Miss rate = **100%**
 - În afara cazului în care vectorul este foarte mic

Initializarea Matricelor in C

- C/C++ utilizează Row-Major
- Prima implementare (i, j)

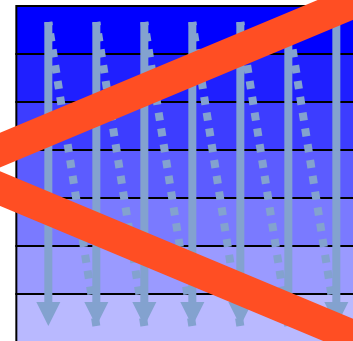
```
int a[100][100];  
for (i=0;i<100;i++)  
    for (j=0;j<100;j++)  
        a[i][j] = 11;
```



Buna Localitate
A Datelor

- A doua implementare (j, i)

```
int a[100][100];  
for (j=0;j<100;j++)  
    for (i=0;i<100;i++)  
        a[i][j] = 11;
```



Masurarea Performantelor

- C/C++ utilizează Row-Major

- Prima implementare

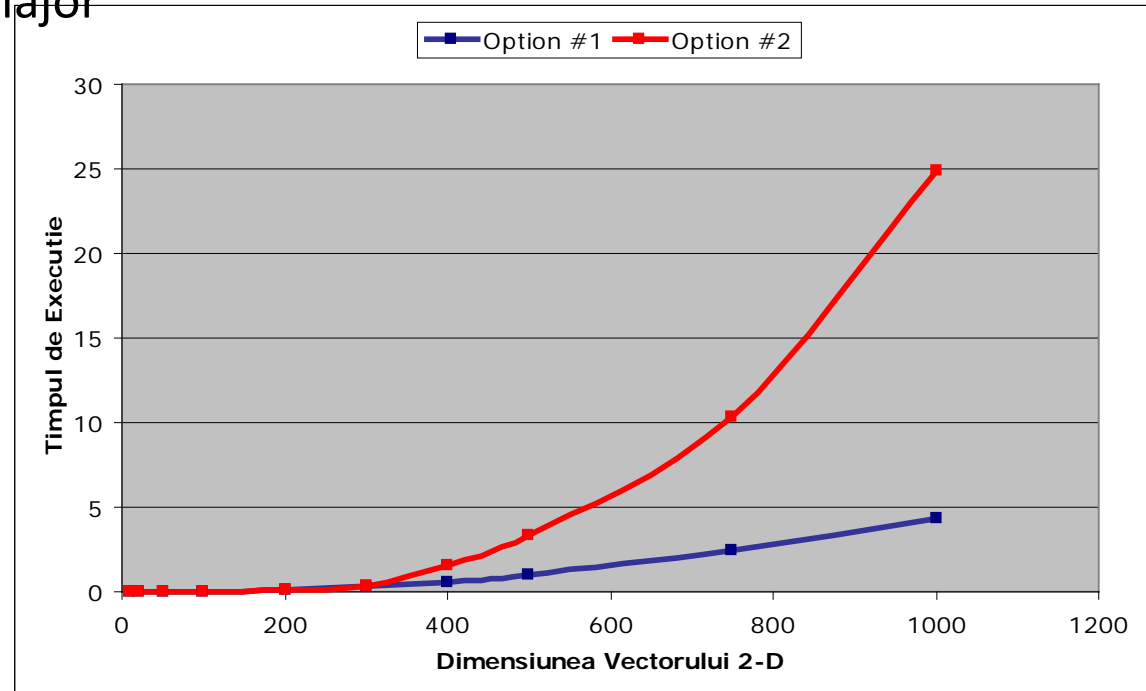
```
int a[100][100];  
for (i=0;i<100;i++)  
    for (j=0;j<100;j++)  
        a[i][j] = 11;
```

- A doua implementare

```
int a[100][100];  
for (j=0;j<100;j++)  
    for (i=0;i<100;i++)
```

```
        a[i][j] = 11;
```

- Alte limbaje de programare utilizează column major
- “
- FORTRAN de exemplu



Experimente pe un PC normal

What About Dynamic Arrays?

- In unele limbaje, putem declara vectori cu dimensiuni variabile
 - FORTRAN:

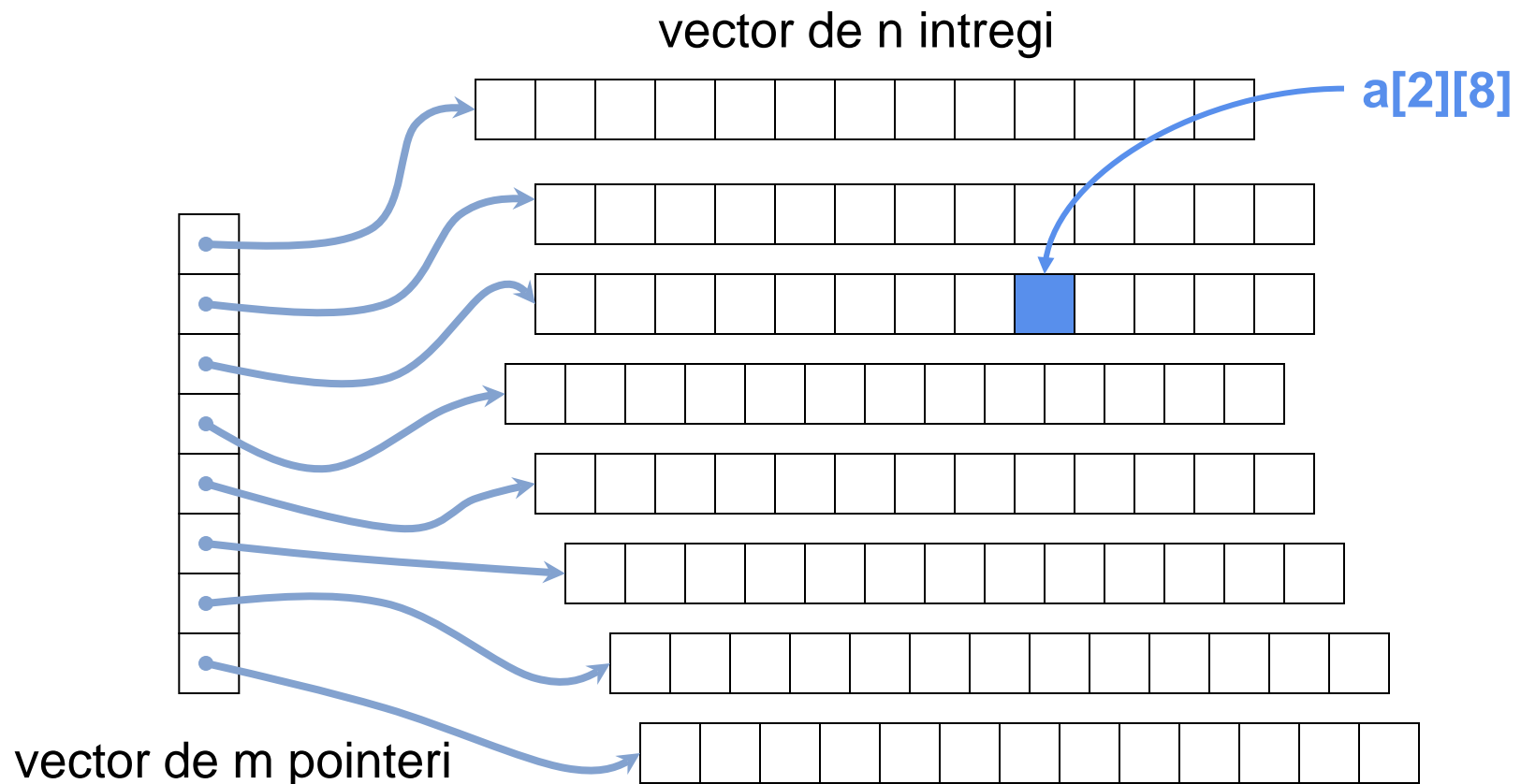
```
INTEGER A(M,N)
```

- C-ul de exemplu nu permite acest lucru
- In C, trebuie sa alocam explicit memoria ca un **vector de vectori**:

```
int **a;  
a = (int **)malloc(m*sizeof(int));  
for (i=0;i<m;i++)  
    a[i] = (int *)malloc(n*sizeof(int));
```

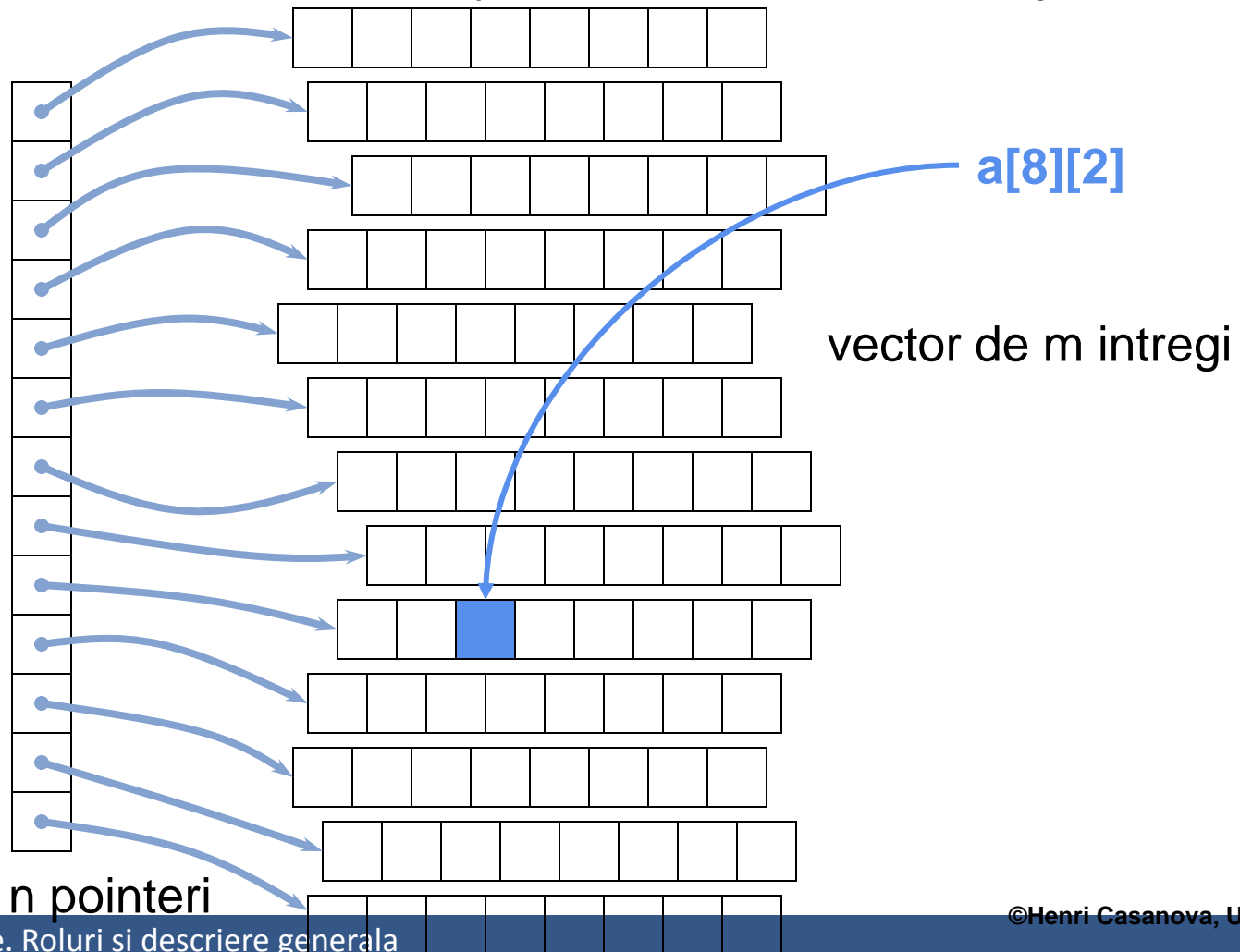
Asezarea Memoriei

- O soluție “non contiguous” de tip row major



Asezarea Memoriei

- Se poate înșă face și într-o implementare column-major

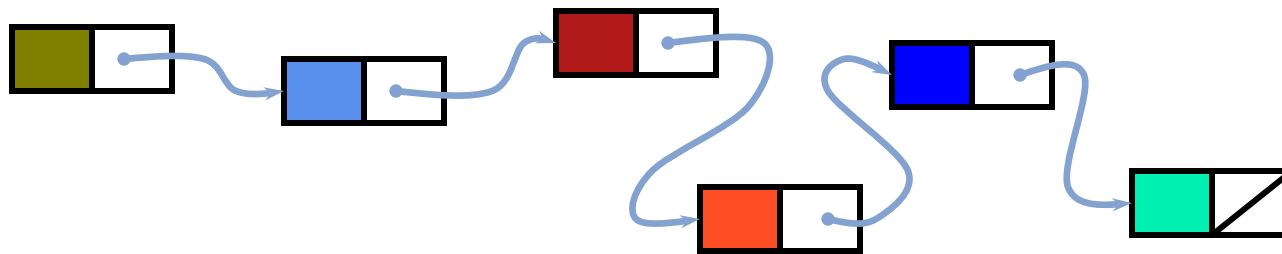


Vectori Dinamici

- Programatorul trebuie sa aleaga
 - Asezarea row-major sau column-major
 - Ordinea in care face initializarea vectorilor
- In Java, toti vectorii sunt alocati dinamic
- Vectorii de dimensiuni mari
 - **Dinamici** – de dimensiuni N-D sunt vectori (N-1)-D de vectori 1-D
 - **Statici** – o generalizare a solutiei row/column-major

Alte Exemple: Liste Inlantuite

- Sa luam in considerare exemplul unei liste inlantuite



- Intr-o implementare tipica, fiecare element al listei va fi alocat dinamic la momentul inserarii
- Elementele listei **nu** vor fi contigue
- Parcurgerea listei in ordine va genera in mod normal cate un cache-miss pentru **fiecare element!**
- Acest lucru poate pune probleme majore de performanta daca lista este suficient de lunga si ea este parcursa des...

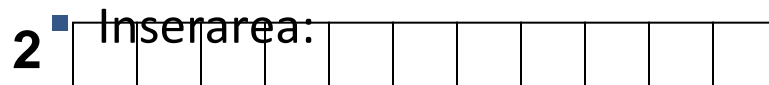
Liste Implementate ca Vectori

- Cum putem avea o localitate mai buna a datelor in liste?
- Le implementam ca vectori
- Tipul de date lista, poate fi implementat sub forma unui vector uni-dimensional
- Sunt trei operatii fundamentale cu liste:
 - Insert (list, current, next)
 - Remove (list, current)
 - Next (list, current)

Liste – Inserarea

- Vom construi o lista de intregi:
- Tipul de date “lista” arata:

▪ `int *array`: un vector de intregi



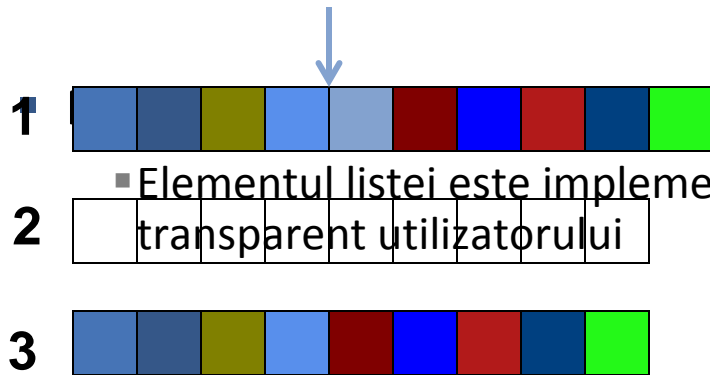
1. **Avem 10 elemente si o inserare**
2. **Alocam un nou vector de 11 elemente**
3. **Copiem elementele vechiului vector**
4. **Copiem noul element**
5. **Eliberam vechiul vector**
6. **Facem lista sa “pointeze” catre noul vector**

■ Stergerea;

Platformă de e-learning și curricula e-conținut pentru învățământul superior tehnic



Liste – Stergerea & Next



1. Avem o lista de 10 elemente și o stergere
2. Alocăm un nou vector de 9 elemente
3. Copiem elementele din vechiul vector
4. Eliberam vechiul vector
5. Facem lista să "pointeze" către noul vector

Liste – Concluzii

- O optimizare simpla:
 - In C, puteti utiliza “realloc” pentru a minimiza copierea elementelor
- Trade-off
 - Inserarea poate dura considerabil mai mult decat in implementarea traditionala
 - Au loc mai multe copieri de date
 - Lucrurile stau la fel si la stergere
 - Operatia Next este insa aproape instantanee
- Utilizatorul trebuie astfel sa identifice “the common case” si sa selecteze implementarea corespunzatoare
- Alte optimizari posibile
 - La inserare: dublati dimensiunea vectorului daca acesta e “prea mic”
 - La stergere: permiteti utilizarea unui vector “mai mare”

Ierarhia de Memorii vs. Programatori



- Toate lucrurile prezentate aici pot fi foarte frumoase (poate)...
- Dar de ce ar fi importante pentru programatori?
- Am văzut că putem determina ordinea de execuție a buclelor
 - Poate că acest lucru îl poate face compilatorul pentru noi
- Totuși, un programator care dorește performanțe de la codul său, **trebuie** să știe cum arată ierarhia de memorii pe care programează!
 - Dacă știm dimensiunea cache-ului L2 (256KB), putem descompune problema în subprobleme de această dimensiune pentru a exploata localitatea datelor
 - Dacă știm că o linie de cache are 32 bytes, putem calcula precis numărul de cache-miss-uri cu o formulă și astfel setăm un parametru optim pentru programul nostru
- Pentru a avea o experiență activă cu ierarhia cache-urilor, e util să vedem cum putem scrie programe care să masoare în mod automat aceste caracteristici



- Sa presupunem urmatorul fragment de cod

```
char a[N];  
for (i=0;i<N;i++)  
    a[i]++;
```

- Presupunem ca L este dimensiunea liniei de cache in bytes

- Numaram numarul de cache-miss-uri:

- a[0]: miss (incarcam o noua linie de cache)
- a[1]: hit (in linie de cache)
- ...
- a[L-1]: hit (in linie de cache)
- a[L]: miss (incarcam o noua linie de cache)
- ...

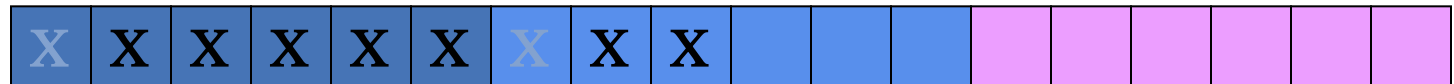
- Numarul de miss-uri este astfel: $\sim N / L$

Masurarea L-ului?

- Codul anterior accesează elementele vectorului cu **pasul 1**
 - **Pasul** este diferența în bytes între două adrese accesate în două accesuri succesive la memorie
- Ce se întâmplă când avem un pas 2?

```
char a[N];  
for (i=0; i<N; i+=2)  
    a[i]++;
```
- Avem practic un număr dublu de *miss-uri* față de cazul anterior!
 - O modificare minoră (aparent) are un impact major asupra performanțelor

Pas 1

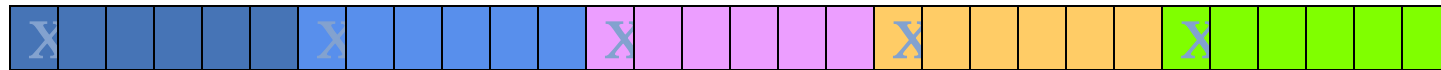


Pas 2

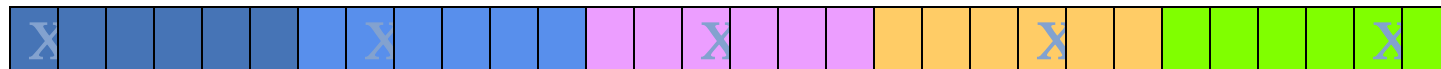


Masurarea L-ului?

- Putem astfel crește pasul până când...
- Pasul ajunge să fie egal cu L:
 - Fiecare acces are nevoie de o linie proprie de cache
 - N accesuri la memorie → N cache-miss-uri



- Ce se întâmplă dacă pasul este L+1?
 - Fiecare acces are nevoie de o linie proprie de cache
 - N accesuri la memorie → N cache-miss-uri



Masurarea L-ului?

- Cea mai buna performanta: $pas=1$
- Cea mai proasta performanta: $pas \geq L$
- Daca masuram performanta codului pentru diverse valori ale pasului, obtinem un grafic de genul:



- Gaseste inceputul platoului (pasul in bytes) pentru care performanta codului nu se mai inrautateste odata cu cresterea pasului
- Aces pas (in bytes) este dimensiunea liniei de cache!

Masurarea L-ului?

- Cum scriem un program pentru măsurarea performanței pentru mai multe valori pentru pas?
- Performanța – timpul mediu pentru accesul la memorie
- Alocati vectori **mari** de caractere
- Creați bucle cu valori pentru pas între 1 și 256
- Pentru fiecare pas ales, parcurgeți în mod repetat vectorul
 - Faceți operații cu fiecare element al vectorului (etc)
 - Măsurați timpul cât durează aceste operații
 - Măsurați câte operații ati operat în total
 - Impartiti numărul de operații la timpul măsurat

Masurarea Dimensiunii Cache-ului

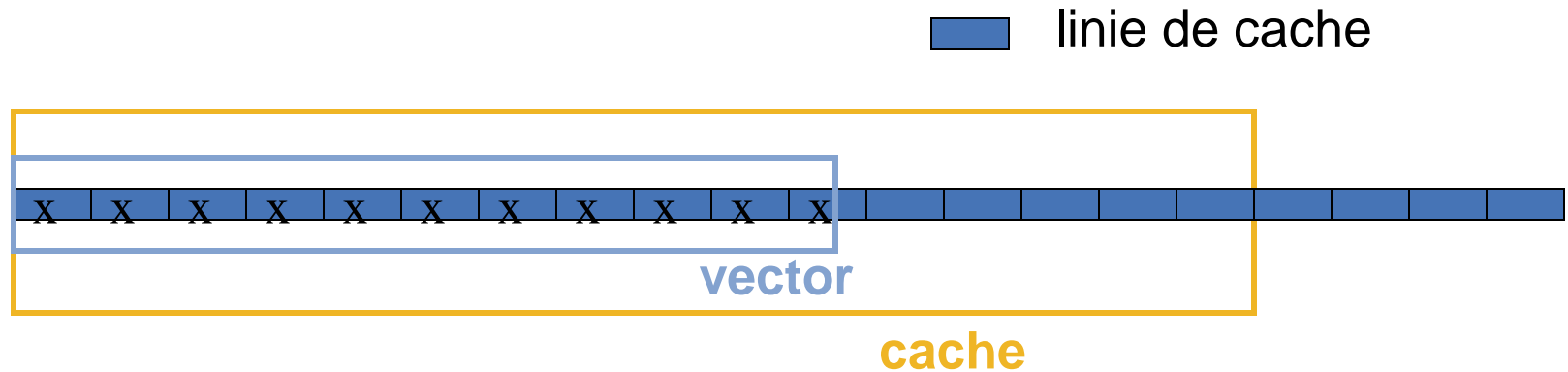
- Dacă L este dimensiunea liniei de cache
- Sa consideram urmatorul cod

```
char x[1024];  
  
for (step=0;step<1000;step++)  
    for (i=0;i<1024;i+=L)  
        x[i]++;
```

- Dacă cache-ul este **mai mare de 1024 de bytes**, după prima iteratie a buclei “**step**”, tot vectorul x e în cache și **nu vom mai avea** miss-uri deloc
- Dacă cache-ul este **mai mic decât 1024 de bytes**, vom avea mereu un număr de miss-uri

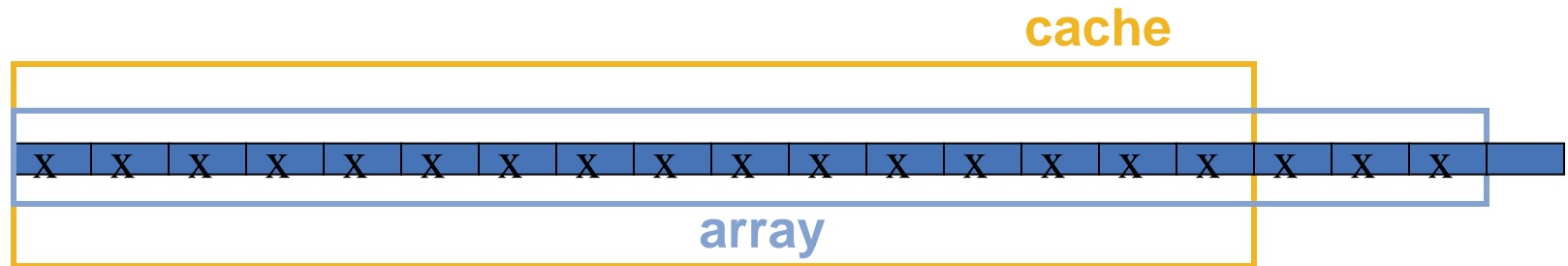
Masurarea Dimensiunii Cache-ului

- Exemplificare:



- Cache-ul este suficient de mare pentru a contine 16 linii
- Vectorul intra in 11 linii de cache
- Astfel, vor fi 11 miss-uri si apoi vor exista doar hit-uri
- Pentru un numar mare de iteratii, hit-rate-ul va fi aproape de 100%

Masurarea Dimensiunii Cache-ului



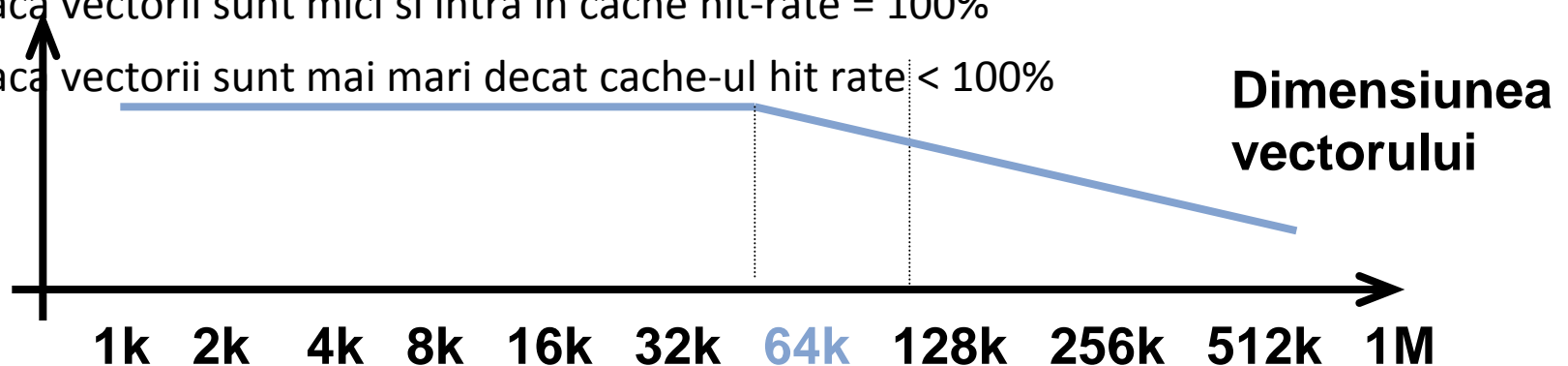
- Cache-ul poate contine 16 linii
- Vectorul intra in sa doar in 19 linii de cache
- Initial vor fi 19 miss-uri
- Apoi va trebui sa citim aceleasi 19 linii de cache
- Doar 16 intra in sa in cache
- Vom avea astfel 13 hit-uri si 3 miss-uri
- Acest comportament va fi repetat pentru fiecare iteratie
- Pentru un numar mare de iteratii, se ajunge la un hit-rate de aproximativ 81%

Masurarea Dimensiunii Cache-ului

- Astfel:

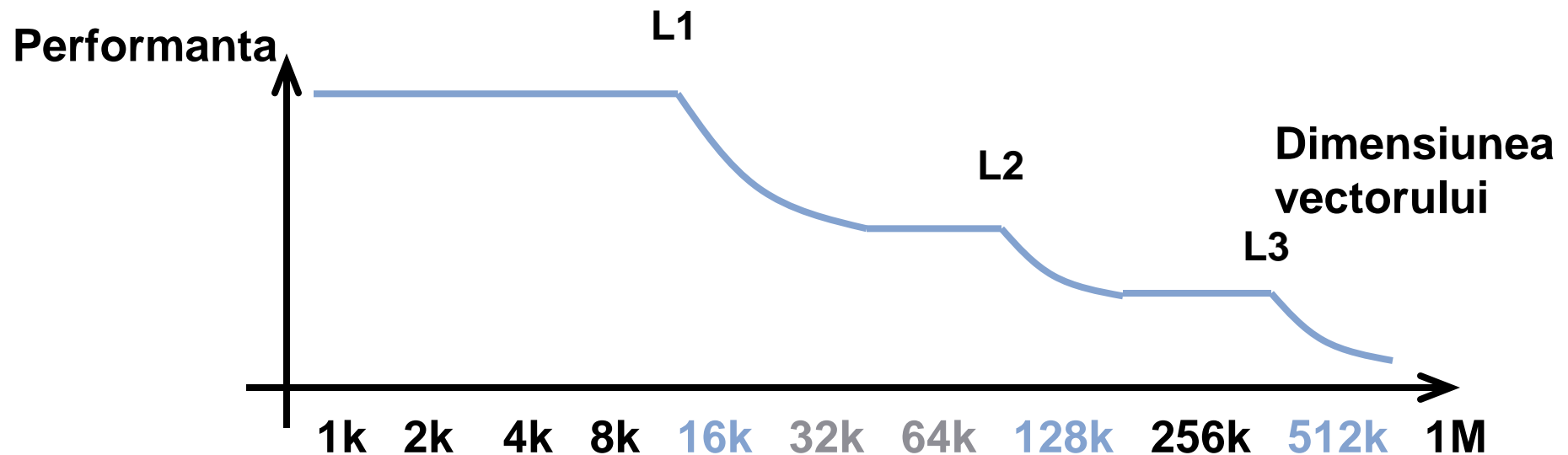
Performanta

- Dacă vectorii sunt mici și intra în cache hit-rate = 100%
- Dacă vectorii sunt mai mari decât cache-ul hit rate < 100%



Mai multe nivele de Cache?

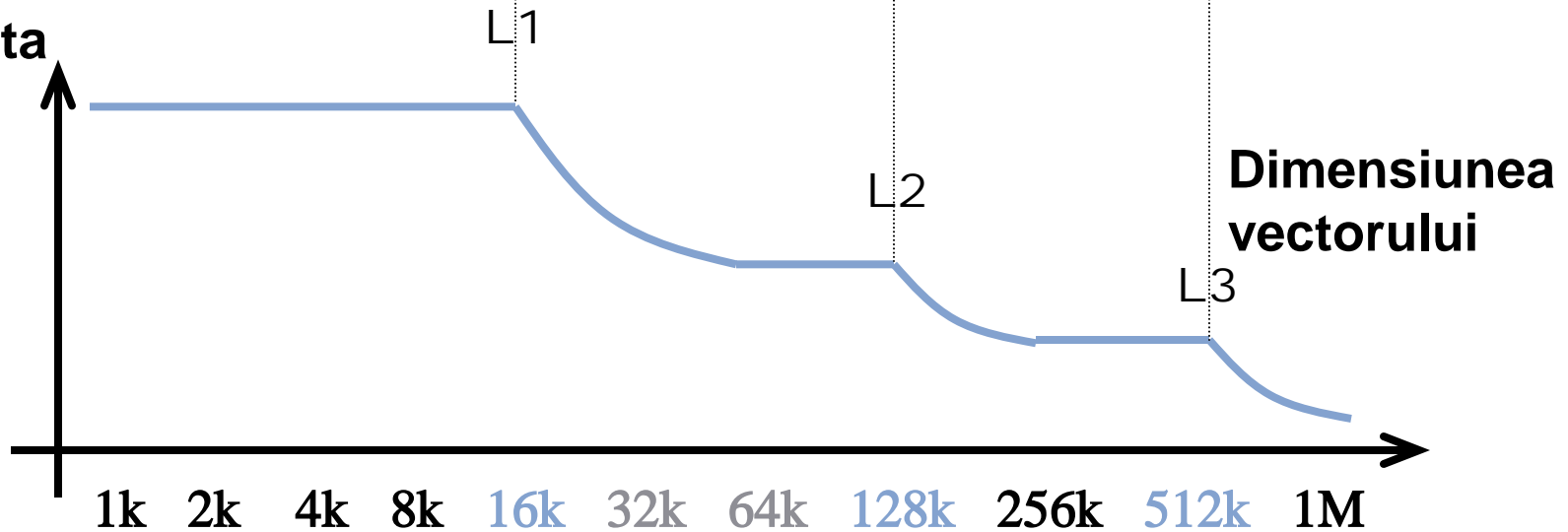
- In general însă, există mai multe nivele de cache: L1, L2, L3
- Având în vedere acest fapt, graficul anterior devine



Mai multe nivele de Cache?

Totul intra în L1	Doar o parte intra în L1	
Totul intra în L2		Doar o parte intra în L2
Totul intra în L3		Doar o parte in

Performanța



Cache-uri in Procesoare Reale

- Exista 2 sau 3 nivele de cache
- Cache-urile aproape de procesor sunt in general mapate direct, si cele mai departate sunt asociative
- Cache-uri diferite de date/instructiuni aproape de procesor, si unificate in rest
- Write-through si write-back sunt la fel de des intalnite, dar nu va exista niciodata o implementare write-through pana la memoria principala
- Liniile de cache aveau in mod normal 32-byte, dar acum exista foarte des linii de 64- si 128-bytes
- Cache-uri neblocante
 - La un miss, nu bloca procesorul ci executa instructiuni de dupa load/store-ul curent
 - Blocheaza procesorul doar daca aceste operatii utilizeaza date din load
 - Cache-urile trebuie sa poata sa serveasca multiple accese “invechite”

Performanța Cache-urilor

- Performanța cache-ului este data de accesul mediu la memorie
 - Programatorii sunt interesați în general doar de timpul de execuție – în timp ce timpul de acces la memorie este o componentă extrem de importantă
- Formula e simplă:
 - $\text{timpul de acces la memorie} = \text{hit time} + \text{miss rate} * \text{miss penalty}$
 - Miss-rate este procentul de miss-uri pe acces la date și **NU** la instrucțiuni!
- La fel ca în cazul timpului de execuție procesor, termenii ecuației **NU** sunt independenți
 - Nu putem spune: reduc numărul de instrucțiuni, fără a crește numărul de instrucțiuni pe ciclu sau viteza ceasului sistemului
 - Similar, nu putem spune: reduc miss-rate-ul fără a crește hit-time

Impactul Miss-urilor

- Miss penalty depinde de tehnologia de implementare a memoriei
- Procesorul masoara numarul de cicluri pierdute
- Un procesor mai rapid e “lovit” mai tare de **memorii lente** si **miss-rate-uri crescute**
- Astfel, cand incercati sa estimati performantele unui sistem de calcul, comportamentul cache-urile **trebuie** luat in considerare (Amdahl – CPU vs. memorie)
- **De ce ne intereseaza?**
- Pentru ca putem rescrie/rearanja codul pentru a utiliza mai bine localitatea datelor

De ce ne Intereseaza?

- Putem citi dimensiunile cache-urilor din specificatiile masini pe care rulam
 - Dar... ceva poate lipsi
 - Se poate sa nu avem acces la specificatii
- E util sa avem o optimizare automatizata
 - Ce trebuie sa fac pentru a scrie un program ce sa ia in considerare Cache-ul?
 - Utilizeaza constanta `CACHE_SIZE` pentru a seta dimensiunea cache-ului sistemului
 - Utilizeaza aceasta constanta cand aloca memoria

```
char x[CACHE_SIZE]
for (i=0;i < CACHE_SIZE; i++)
```
 - Inainte de a compila un program, ruleaza un alt utilitar pentru a descoperi dimensiunea cache-ului
 - Seteaza apoi constanta `CACHE_SIZE` la valoarea astfel determinata

```
#define CACHE_SIZE 1024*32
```
- Un programator (bun) nu poate ignora cache-ul
- Desi unele structuri de date par naturale, ele pot fi extrem de ineficiente in cache (liste, pointeri, etc) si de aceea ele ar trebui evitate cand se doreste o performanta crescuta a codului