

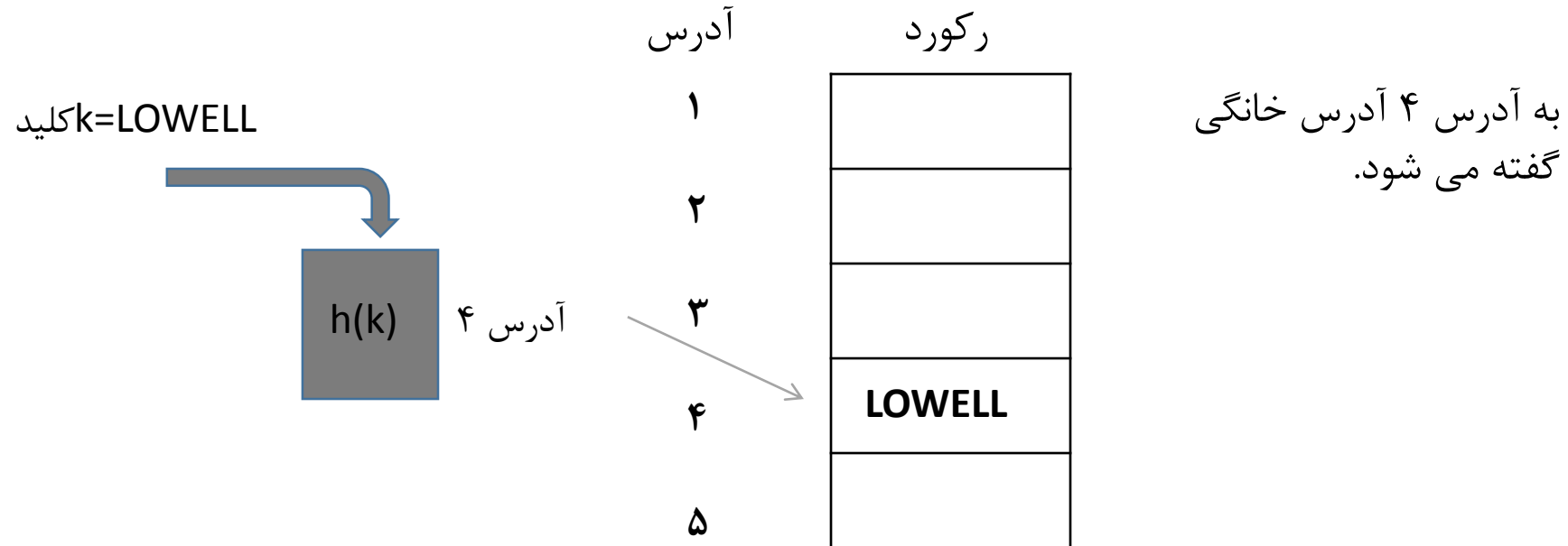
درهم سازی

سید مهدی وحیدی پور

ارایه سوم: درهم سازی، روشها و کاربردها

درهم سازی

□ درهم سازی تابع $h(k)$ است که کلید k را به یک آدرس انتقال می دهد. از این آدرس به عنوان مبنایی برای ذخیره و بازیابی رکوردها استفاده می شود.



تفاوت درهم سازی و اندیس سازی

- ❑ آدرس‌های به دست آمده از درهم سازی تصادفی هستند.
- ❑ با درهم سازی دو کلید مختلف ممکن است به یک آدرس انتقال داده شوند. که در آن صورت برخورد (Collision) رخ می‌دهد.

چندین روش مختلف برای کاهش تعداد برخوردها

➤ پراکنده کردن رکوردها

بتوانیم الگوریتمی پیدا کنیم که رکوردها را به صورت تصادفی و پراکنده بین آدرس ها توزیع کند که تعداد زیادی از رکوردها در اطراف یک آدرس مشخص تجمع نکند.

➤ استفاده از حافظه اضافی

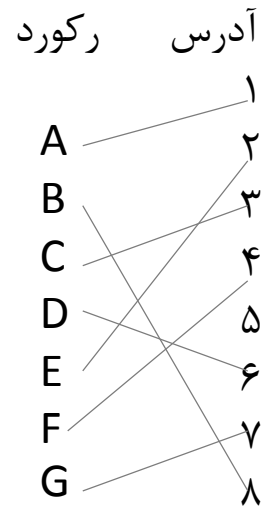
در مقابل تعداد رکورد از حافظه بیشتری کمک می گیریم که پیدا کردن الگوریتم در هم سازی ساده تر باشد.

➤ قرار دادن بیشتر از یک رکورد در یک آدرس

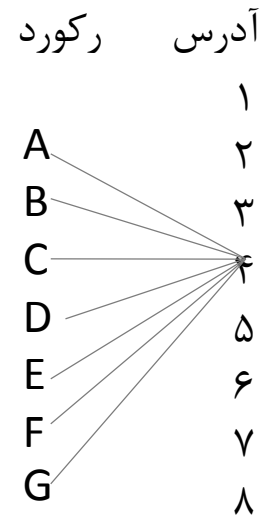
فایل به اندازه کافی بزرگ باشد تا چندین رکورد را در خود ذخیره کند. به آدرس هایی که چندین رکورد را نگهداری می کنند باکت گفته می شود.

توزیع رکوردها در بین آدرس‌ها

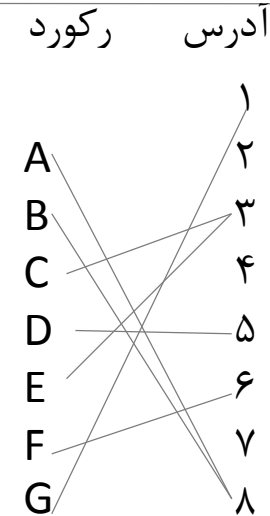
بهترین



بدترین



قابل قبول



روش های درهم سازی

- جستجو در کلید برای یافتن یک الگو: اگر بعضی از قسمت های یک کلید، یک الگو واقعی قابل استفاده را نشان دهد. می توان از تابع درهم سازی استفاده کرد که آن بخش از کلید را استخراج کند.
- تا کردن قسمت هایی از کلید: عددهایی را از قسمت هایی از کلید استخراج می کند و سپس به آنها را به هم اضافه می کند. این روش الگوهای اولیه را خراب می کند.
- تقسیم یک کلید بر یک عدد: اندازه آدرس را بر یک عدد تقسیم می کنیم و از باقیماندهی آن استفاده می کنیم.
- مجذور کردن کلید و گرفتن عدد میانی
- تبدیل مبنا

پیش بینی توزیع رکوردها

N تعداد ادرس های موجود

R تعداد رکوردهایی که باید ذخیره شوند.

X تعداد رکوردهایی که به یک ادرس مشخص اختصاص داده شده اند.

$$P(x) = \frac{(r/n)^x e^{-(r/n)}}{X!} \quad \text{توزیع پواسون}$$

$$P(0) = \frac{1^0 e^{-1}}{0!} = 0.368$$

احتمال اینکه هیچ کلیدی در یک ادرس مشخص در هم سازی نشود:

$$P(1) = \frac{1^1 e^{-1}}{1!} = 0.368$$

احتمال اینکه دقیقا یک، دو یا سه کلید در یک ادرس مشخص در هم سازی شوند:

$$P(2) = \frac{1^2 e^{-1}}{2!} = 0.184$$

$$P(3) = \frac{1^3 e^{-1}}{3!} = 0.061$$

دانسیتة فشردگی

$$\text{دانسیتة فشردگی} = R/N = \frac{\text{تعداد رکورد ها}}{\text{تعداد فضا}}$$

اندازه فضای استفاده شده در فایل را نشان می دهد و تنها راه تخمین کارایی محیط درهم سازی است.

دانسیتة فشردگی - ادامه

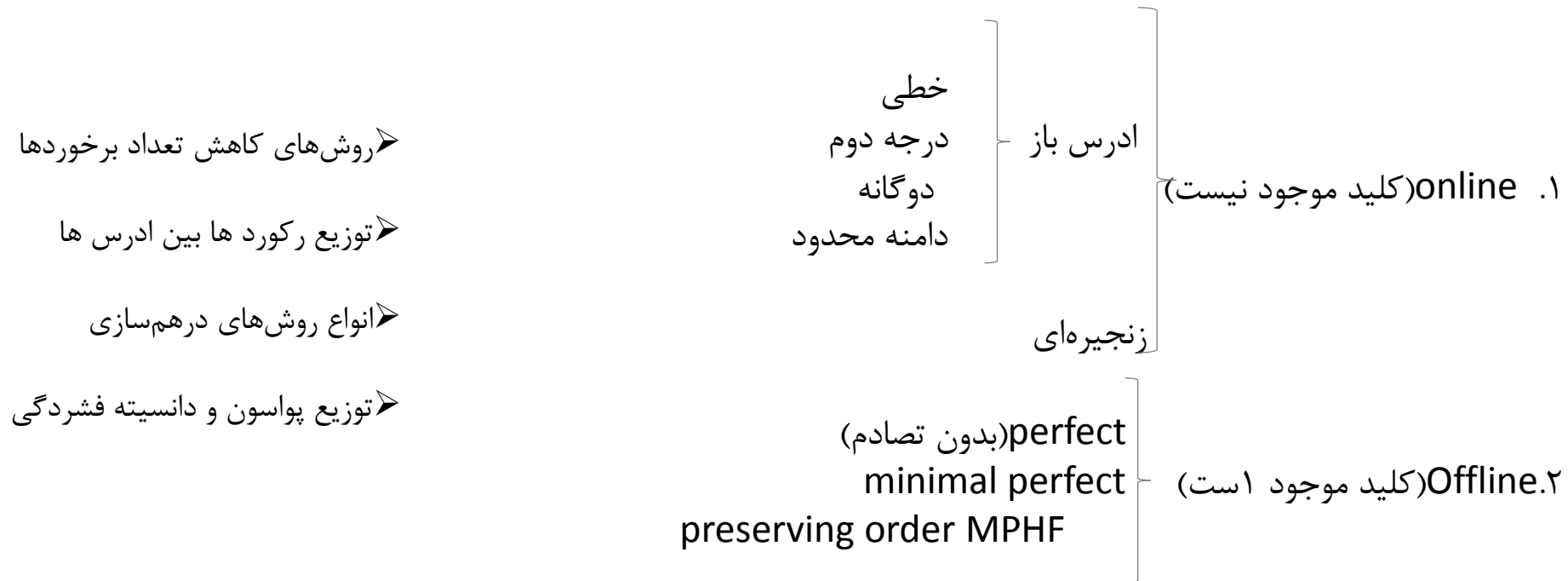
اگر دانسیته فشردگی ۵۰ درصد باشد و هر ادرس تنها قادر باشد یک رکورد را در خود ذخیره کند انتظار داریم که ۲۱ درصد تعداد کل رکوردها در جایی غیر از ادرس خانگی خودشان ذخیره شوند.

مترادف ها به صورت درصدی از رکوردها	دانسیته فشردگی (درصد)
۴.۸	۱۰
۱۳.۶	۳۰
۲۱.۴	۵۰
۲۸.۱	۷۰
۳۱.۲	۸۰
۳۴.۱	۹۰
۳۶.۸	۱۰۰

اثر دانسیته فشردگی نسبت به رکورد هایی که در ادرس خانگی خود ذخیره نشده اند:

از تابع پواسون و دانسیته نتیجه می شود که برخورد اجتناب ناپذیر است در نتیجه باید به روش های رفع برخورد بیاندیشیم.

روش های رفع برخورد



- ذخیره کردن بیشتر از یک رکورد در هر ادرس (باکت)
- پیوند با ناحیه سرریز
- جدول های پراکندگی: اندیس سازی

درهم سازی آدرس دهی باز

درهم سازی خطی

فرمول کلی

$$H_i(k,i) = (h(k) + c*i) \bmod m, i = 0,1,\dots$$

k مقدار کلید

h(k) یک تابع درهم ساز ساده

i نشان دهنده تعداد تکرار تولید آدرس برای کلید k

مشکلات:

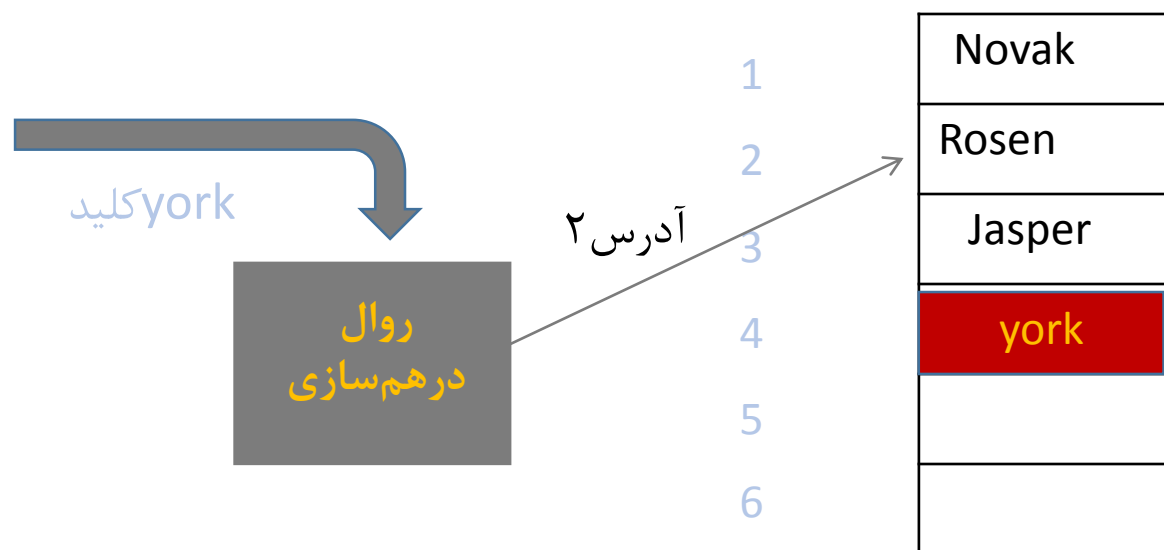
۱_خوشه بندی اولیه: به ازای مقدار ثابت C و تمام مقادیر اولیه برای تابع $h(k,i)=h(k)$ یک دنباله ادرس مشخص بر اساس میزان جابجایی C تولید خواهد نمود.

۲_خوشه بندی ثانویه: به ازای دو کلید مختلف k_1 و k_2 رابطه $h(k_1)=h(k_2)$ برقرار باشد در این صورت تمام ادرس های دو دنباله یکسان خواهد بود.

درهم سازی آدرس دهی باز

روش سرریز فزاینده

- درهم سازی خطی با $c=1$
- اگر به یک رکورد آدرسی تخصیص داده شود که قبلاً توسط رکوردی اشغال شده است آدرس‌های بعدی مورد جستجو قرار می‌گیرند تا اولین آدرس خالی پیدا شود.



□ هنگام جستجو کلید `york` جستجو از آدرس ۲ شروع می‌شود زیرا `york` هنوز در آدرس ۲ درهم‌سازی می‌شود.

روش‌های مختلف جانشین‌سازی

FCFS ➤: کلید قدیمی در سلول نگه داشته میشود و کلید جدید سلول‌های بعدی را بررسی می‌کند.

LCFS ➤: کلید جدید وارد سلول میشود و کلید قدیمی سلول‌های دیگر را بررسی می‌کند.

Robin hood ➤: در بین دو کلید کلیدی که بزرگترین تغییر مکان را دارد در سلول نگه داشته می‌شود و کلید دیگر سلول‌های دیگر را بررسی میکند.

درهم سازی آدرس دهی باز

درهم سازی دوگانه

فرمول کلی

$$HD(k,i) = (g(k) + ih(k)) \bmod m$$

k مقدار کلید

g(k) و h(k) توابع درهم ساز ساده

i نشان دهنده تعداد تکرار تولید آدرس برای کلید k

- رکورد های سرریز در آدرس دورتری از آدرس خانگی قرار دارند.
- هنگام ایجاد برخورد تابع درهم سازی دیگری روی کلید عمل میکند.

مشکلات:

□ محلی بودن را نقض می کند

□ تولید آدرس های یکسان و تکراری در دنباله های آدرس.

درهم سازی آدرس دهی باز

درهم سازی دامنه محدود

به منظور استفاده از روش دامنه محدود، روند زیر برای بدست آوردن آدرس کلید K در فضایی با $m = p^n$ مکان اجرا می شود. P عدد اول و $n \geq 1$.

• به ازای هر کلید k دو عدد a, b را محاسبه کنید:

$$a = k \bmod m, \quad b = k \bmod (m-2)$$

• واضح است که روابط $0 \leq b \leq m-2$, $0 \leq a \leq m-1$ برقرار می شود.

• مجموعه S را یک مجموعه دامنه محدود $GF(P^n)$ می نامند. هر چند جمله ای از محدوده $GF(P^n)$ را با n عدد ضریب ان نشان می دهند. $a = (a_{n-1}, a_{n-2}, \dots, a_0)_P$. لذا یک عدد N رقمی در مبنای P نشان دهنده یک چند جمله ای در محدوده $GF(P^n)$ است.

• اعداد a و b در مبنای p بر اساس محدوده $GF(p^n)$ نشان دهنده دو چند جمله ای به نامهای $g_k(X)$ و $h_k(X)$ می باشند.

$$a = (a_{n-1}, a_{n-2}, \dots, a_0)_P$$

$$b = (b_{n-1}, b_{n-2}, \dots, b_0)_P$$

$$g_k(X) = a_{n-1}X^{n-1} + a_{n-2}X^{n-2} + \dots + a_0$$

$$h_k(X) = b_{n-1}X^{n-1} + b_{n-2}X^{n-2} + \dots + b_0$$

درهم سازی آدرس دهی باز

درهم سازی دامنه محدود

- شماره دنباله آدرس ها، که نشان دهنده ی چندمین آدرس تولیدی است، را نیز در مبنای p نمایش داده و چندجمله ای $f_i(x)$ در فضای $GF(p^n)$ تولید می شود.

$$i = (i_{n-1}, i_{n-2}, \dots, i_0)_p$$

$$f_i(X) = i_{n-1}X^{n-1} + i_{n-2}X^{n-2} + \dots + i_0$$

- حال بر اساس قوانین حاکم بر $GF(p^n)$ و با در نظر گرفتن $f_\lambda(X)$ به عنوان چند جمله ای ثابت در فضای $GF(p^n)$ رابطه زیر نشان دهنده تابع درهم ساز جدید است.

$$H_f(k,i) = (g_k(X) + f_\lambda(X) f_i(X) h_k(X)) \text{ mod } t(x)$$

- $H_f(k,i)$ چندجمله ای از فضای $GF(p^n)$ می باشد که ضرایب آن نشان دهنده عددی است در مبنای p که همان آدرس تولید شده تابع درهم ساز است.

- $t(x)$ چندجمله ای تجزیه ناپذیر در مبنای p مانند

چندجمله ای تجزیه ناپذیر

$$x^4 + x + 1, x^4 + x^3 + x^2 + x + 1, x^4 + x^3 + 1$$

$$x^8 + x^7 + x^6 + x^5 + x^4 + x^3 + x^2 + 1$$

ذخیره کردن بیشتر از یک رکورد در یک آدرس: باکت ها

- گاهی ذخیره سازی گروهی از رکوردها در بلوکها بهتر از ذخیره سازی تک تک آنها است.
- در واقع چند رکورد از آدرس مشترک دارند.
- برای جستجو یک رکورد کل باکت در حافظه قرار می گیرد و رکوردهای موجود در باکت جستجو می شوند.
- حالت سرریز در باکت وجود دارد اما احتمال وقوع آن بسیار کمتر است.

حذف رکورد

حذف فایل پیچیده تر از اضافه کردن آن است، زیرا:

- محل رکورد حذف شده نباید مانع جستجوهای بعدی شود.
- امکان استفاده از فضای آزاد شده باید وجود داشته باشد.

راه حل: استفاده از علائم ویژه

Adams
Jones
#####
Smith

سرریز فزاینده زنجیره‌ای

- رکورد های مترادف در یک لیست پیوندی قرار می گیرند. یعنی هر آدرس خانگی حاوی عددی است که محل رکورد بعدی که در این آدرس خانگی قرار دارد مشخص می کند.

آدرس خانگی	آدرس واقعی	داده	آدرس مترادف بعدی	طول جستجو
۲۰	۲۰	Adams	22	1
۲۱	۲۱	Bates	23	1
۲۰	۲۲	Cole	25	2
۲۱	۲۳	Dean	-1	2
۲۰	۲۴	Evans	-1	1
۲۴	۲۵	Flint	-1	3

متوسط طول جستجو: ۱.۷

مزیت: در هر جستجو فقط رکوردهایی که دارای کلیدهای مترادف هستند دستیابی می شوند.

پیوند با ناحیه سرریز

تمام رکورد های سرریز را به یک ناحیه سرریز منتقل می کنیم.

مجموعه آدرس های خانگی ناحیه داده اصلی و مجموعه آدرس های سرریز ناحیه سرریز نام دارد.

آدرس خانگی

ناحیه داده اصلی

۲۰

Adams

0

→

0

۲۱

Bates

1

→

1

۲۲

2

۲۳

3

۲۴

Evans

-1

ناحیه سرریز

Cole

2

Dean

-1

Flint

-1

جدول های پراکندگی: اندیس سازی

- فایل در هم سازی حاوی هیچ رکوردی نیست و فقط حاوی اشاره گرهایی به رکوردهاست.
- این فایل مانند یک اندیس است که توسط درهم سازی جستجو می شود.

