



الڱور يټم
عضويت

محسن هوشمند
دانشكده علم رايانه و تكنولوژي اطلاعات
دانشگاه تحصيلات تكميلي علوم پايه زنجان

عضویت

مسئله عضویت برای مجموعه داده،

- تصمیم‌گیری در مورد تعلق یا عدم تعلق مقداری به آن مجموعه

برای مجموعه‌های کوچک،

- حل مسئله با جستجوی مستقیم و مقایسه مقدار داده شده با هر عضو
- وابستگی چنین رویکرد ساده‌ای به تعداد عناصر مجموعه و به طور متوسط نیاز به $O(\log n)$ مقایسه (روی داده‌های مرتب

برای مجموعه‌های عظیم از مقادیر

- ناکارآمدی رویکرد مذکور
- نیاز به زمان و حافظه $O(n)$ برای ذخیره عناصر

راه‌حل‌های ممکن

- مانند تقسیم چنین مجموعه‌هایی به قطعات و اجرای مقایسه‌ها به صورت موازی
- کمک کننده در کاهش زمان محاسبه
- عدم قابل اجرا بودن چنین روش‌هایی
- ناممکنی ذخیره چنین مجموعه‌های عظیمی از عناصر
- عدم نیاز به یافتن اصل تطابق و کفایت صرف وجود تطابق \Leftarrow امکان ذخیره امضای عناصر به جای مقدار کامل آنها

عضویت - مثال - مسئله مرور ایمن

ایجاد مرورگر وب

- برخی از URL-ها حاوی بدافزار هستند

در پی هشدار به کاربر در صورت درخواست وی برای بازدید از آن صفحات

- یا حتی از جلوگیری بازدیدش

راه حل فوری با حداقل سازی ترافیک شبکه

- ذخیره همه این URL-ها در برنامه و پس از وارد کردن URL توسط کاربر، صرفا بررسی بدافزاری آن

هدایت در صورت عدم بدافزاری

امکان چنین پیاده سازی ساده ای تا زمان کم بودن تعداد URL-های مخرب کم

- اما، اینصداق نبودن در برنامه های دنیای واقعی

نیاز به ساختار ویژه ذخیره URL-های مخرب (یا ایده آل، فقط برخی اطلاعات درباره آنها) بدون رشد خطی در اندازه

سایر الزامات شامل

- پشتیبانی از بررسی فهرست شدن URL و سرعت تا حد امکان اجتناب از انتشار طولانی مدت کاربران مدت طولانی

- کاربردهای مسئله عضویت نقش اساسی در علوم کامپیوتر و شاخه های مختلف

عضویت-مثال - توالی‌های دی‌ان‌ای

طبقه‌بندی توالی‌ها به عنوان «جدید» یا متعلق به ژنوم شناخته شده

- از مسائل مهم در مطالعات متاژنومیک

- فیلتر کردن داده‌هایی دیده شده

مرحله پیش‌پردازش راه‌انداز عضویت

- در صورت اجرای کارآمد منجر به کاهش پیچیدگی داده‌ها قبل از انجام تحلیل دقیق‌تر

حل مسئله جستجوی سریع با درهم‌ساز

- ساده‌ترین راه برای انجام آن

نگاشت هر عضو از مجموعه به جدول درهم با تابع درهم‌ساز

- دارای احتمال کمی از خطا (ناشی از برخوردهای احتمالی درهم‌ساز)

- نیاز به حدود $O(\log n)$ بیت برای هر مقدار درهم‌ساز شده

- در عمل غیرقابل اجرا برای مجموعه‌های داده عظیم

مطالب جاری

معرفی راه‌حل‌های متفاوت از جدول‌های درهم با

- نیاز به فضای کمتر
- جستجوهای سریع‌تر
- احتمالات خطای کمتر

چنین داده‌ساختارهای کارآمد از نظر فضا کمک‌کننده به مدیریت حجم

امکان اجرای پرسش‌های عضویت با عملکرد خوب

فیلتر بلوم، بهبودها، و جانشین‌های متاخر

فیلتر بلوم

از سازوکارهای استاندارد مسئله عضویت در سامانه‌های با مجموعه داده‌های بسیار بزرگ کاربرد گسترده آن‌ها، به‌ویژه در شبکه‌ها و پایگاه داده‌های توزیع شده،

کارآمدی قابل توجه در شرایطی

- نیاز به نوعی کارکرد مشابه جدول درهم
- اما بدون در اختیار داشتن فضای کافی

بارتون بلوم معرفی آن در سال ۱۳۴۹

- اما مهجوری تا این اواخر
- اشتهار طی دو دهه اخیر
- در پی افزایش نیاز به مهار و فشرده‌سازی مجموعه داده‌های عظیم
- شکوفائی! ساختار

جلب توجه جامعه پژوهشی علم رایانه

- طراحی گونه‌های متعدد از آن با هدف برطرف کردن برخی محدودیت‌های این ساختار و سازگار کردن آن با نیازمندی‌های متفاوت

فیلتر بلوم

داده ساختار احتمالاتی فضا-کارآمد برای نمایش مجموعه داده $D = \{x_1, x_2, \dots, x_n\}$ از n عضو

پشتیبانی از دو عملیات افزودن مقدار به مجموعه، و آزمون عضویت در مجموعه

امکان ذخیره کارآمد مجموعه بزرگ با کنار گذاشتن ذخیره دقیق خود مقادیر

ذخیره مجموعه‌ای (تقریباً) منحصر به فرد از بیت‌ها

▪ هر بیت حاصل تعدادی تابع درهم‌ساز اعمالی الگوریتم روی مقدار

نمایش با آرایه‌ای بیتی با طول m و تعداد توابع درهم‌ساز مختلف $\{h_i\}_{i=1}^k$

طراحی و انتخاب m متناسب با تعداد عناصر مورد انتظار n و $k \ll m$

فیلتر بلوم

توابع درهم‌ساز h_i مستقل از هم و به طور یکنواخت توزیع شده

- حاصل توزیعی یکنواخت و تصادفی
- توابع درهم‌ساز به مثابه نوعی مولد اعداد تصادفی
- کاهش احتمال تصادم در فیلتر

کاهش شدید فضای ذخیره‌سازی

بدون توجه به تعداد عناصر در داده‌ساختار و اندازه آنها

- با رزرو چند بیت برای هر عنصر،
- نیاز به تعداد ثابتی بیت

فیلتر بلوم

داده ساختار آرایه‌ای بیتی به طول m

- ست که در ابتدا همه بیت‌ها برابر با صفر
- بمعنای خالی بودن فیلتر

درج عنصر x در فیلتر،

- محاسبه مقدار درهم متناظر با هر تابع درهم‌ساز h_i یا $j = h_i(x)$
- یک شدن بیت متناظر j در فیلتر

امکان تنظیم چند بارهٔ بیتی به دلیل تصادم‌های درهم‌ساز

فیلتر بلوم

الگوریتم - افزودن عنصر به فیلتر بلوم

ورودی: عنصر $x \in D$

ورودی: فیلتر بلوم با k تابع درهم‌ساز $\{h_i\}_{i=1}^k$

برای i از ۱ تا k انجام بده

$$j \leftarrow h_i(x)$$

$$BloomFilter[j] \leftarrow 1$$

فیلتر بلوم

درج نام پایتخت‌ها در فیلتر

کپنهاگ

محاسبه مقادیر درهم برای یافتن بیت‌های متناظر در فیلتر

$$h_1(\text{Copenhagen}) = \text{MurmurHash3}(\text{Copenhagen}) \% 10 = 7$$

$$h_2(\text{Copenhagen}) = \text{FNV1a}(\text{Copenhagen}) \% 10 = 3$$

تنظیم بیت‌های ۳ و ۷ در فیلتر

0	1	2	3	4	5	6	7	8	9
0	0	0	1	0	0	0	1	0	0

فیلتر بلوم

امکان بیت‌های متناظر مشترک برای ورودی‌های مختلف

مثال: دوبلین

$$h_1(\text{Dublin}) = \text{MurmurHash3}(\text{Dublin}) \% 10 = 1$$

$$h_2(\text{Dublin}) = \text{FNV1a}(\text{Dublin}) \% 10 = 3$$

موقعیت‌های بیت متناظر آن در فیلتر ۱ و ۳

- تنظیم نبودن بیت ۱ از قبل
- به معنای نبودن دوبلین در فیلتر
- مدخل ۳ به عنوان یکی از بیت‌های متناظر از قبل پر

0	1	2	3	4	5	6	7	8	9
0	1	0	1	0	0	0	1	0	0

فیلتر بلوم

آزمون وجود مقدار x در فیلتر،

محاسبه تمامی k تابع درهم‌ساز $\{h_i(x)\}_{i=1}^k$

بررسی بیت‌ها در موقعیت‌های متناظر

اگر همه بیت‌ها روی یک تنظیم شده باشند، امکان عضویت مقدار x در فیلتر
▪ در غیر این صورت، معدوم وجود قطعی مقدار x در فیلتر

عدم قطعیت در مورد وجود عنصر

ناشی از احتمال موقعیت‌هایی از تنظیم برخی از بیت‌های حاصل از مقادیر قبلی اضافه شده
یا به دلیل تصادم‌های شدید، در مواقع تصادم توابع درهم‌ساز با هم

فیلتر بلوم

الگوریتم-۲- آزمون عنصر در فیلتر بلوم

ورودی: عنصر $x \in D$

ورودی: فیلتر بلوم با k تابع درهم $\{h_i(x)\}_{i=1}^k$

خروجی: غلط در صورت نبودن عنصر و درست در صورت احتمال وجود عنصر

برای i از ۱ تا k انجام بده

$$j \leftarrow h_i(x)$$

اگر $BloomFilter[j] \neq 1$

برگرداندن غلط

برگرداندن درست

فیلتر بلوم

مثال - آزمون عضویت در فیلتر

فیلتر بلوم با دو مقدار نمایه شده، کپنهاگ و دوبلین

0	۱	۲	۳	۴	۵	۶	۷	۸	۹
0	۱	0	۱	0	0	0	۱	0	0

فیلتر بلوم

آزمون وجود عنصر کپنهاگ در فیلتر

محاسبه مقادیر درهم آن

$$h_1(\text{Copenhagen}) = 7$$

$$h_2(\text{Copenhagen}) = 3$$

بررسی بیت‌های متناظر در فیلتر بررسی و برابر یک بودن هر دو

▪ ادعا بر احتمال وجود کپنهاگ در فیلتر

فیلتر بلوم

رُم

محاسبهٔ مقادیر درهم‌ساز آن جهت یافتن بیت‌های متناظر در فیلتر

$$h_1(\text{Rome}) = \text{MurmurHash3}(\text{Rome}) \% 10 = 5$$

$$h_2(\text{Rome}) = \text{FNV1a}(\text{Rome}) \% 10 = 6$$

بررسی بیت‌های ۵ و ۶

- تنظیم نبودن بیت ۵
- عدم عضویت مقدار رُم در فیلتر
- عدم نیاز به بررسی بیت ۶

فیلتر بلوم

مقادیر درهم متناظر شهر برلین

$$h_1(\text{Berlin}) = \text{MurmurHash3}(\text{Berlin}) \% 10 = 1$$

$$h_2(\text{Berlin}) = \text{FNV1a}(\text{Berlin}) \% 10 = 7$$

تنظیم بیت‌های متناظر ۱ و ۷ هر دو در فیلتر

پیروزی نتیجه تابع آزمون و امکان وجود مقدار در فیلتر

اما طبق مثال اطلاع از عدم وارد کردن چنین مقداری

مثبت کاذب

نمونه‌ای از تصادم درهم‌ساز است.

در این مورد خاص تنظیم بیت ۱ با $h_1(\text{Dublin})$ و بیت ۷ با $h_2(\text{Copenhagen})$

فیلتر بلوم

در صورت زمان ثبات محاسبه هر تابع درهم‌ساز $\{h_i(x)\}_{i=1}^k$ صادق در توابع درهم‌ساز پرکاربرد

زمان افزودن مقداری جدید یا آزمون مقداری ثابت $O(1)$

مستقل از طول فیلتر m و تعداد مقایدها افزوده شده به فیلتر

وابستگی عملکرد فیلتر بلوم به توابع درهم‌ساز انتخاب شده

کاهش میزان مثبت کاذب مشاهده شده عملی با تابع درهم‌ساز با یکنواختی خوب

سریع‌تر بودن محاسبه هر تابع درهم‌ساز منجر به کاهش زمان کل هر عملیات

توصیه به اجتناب از توابع درهم‌ساز رمزنگاری

فیلتر بلوم

مثال - جلوگیری از رمزهای عبور به خطر افتاده

صفحه ثبت نام سرویس وبی

در پی اجتناب از انتخاب رمزهای عبور ضعیف و به خطر افتاده کاربران

امکان یافتن صدها میلیون رمز عبور هک شده در وب تاریخ

امکان استفاده در حمله لغت نامه‌ای

- حمله جستجو کامل که با امتحان همه مقادیر از فهرستی از پیش آماده
- تلاش‌های تکراری برای شکست احراز هویت

اطمینان از نبود در فهرست با رمز عبور جدیدی

- عدم امنیت ذخیره رمزهای عبور خام
- عدم ذخیره مجموعه داده عظیم دارای رشد خطی با افزودن هر رمز جدید و کندسازی جستجوها

استفاده از فیلتر بلوم فضا-کارآمد

رویداد مثبت کاذب

- تشخیص اشتباه رمز عبور وارد شده نامناسب
- در چنین موارد نادری، درخواست از کاربر برای رمز عبور دیگری
- معمولاً بدون ضرری

فیلتر بلوم

مثال- کاربرد در اپلیکیشن موبایل بیت کوین:

- استفاده از فیلتر بلوم در تبادل داده در شبکه‌های همتابه‌همتا
- سامانه بیت کوین
- از ویژگی‌های مهم بیت کوین، تضمین شفافیت بین کاربران
- تحقق از طریق آگاه‌بودن هر گره از تراکنش‌های سایر گره‌ها
- گره‌هایی روی تلفن‌های هوشمند یا دستگاه‌های مشابه با حافظه و پهنای باند محدود اجرا
- نامناسبی و پرهزینه‌گی نگهداری نسخه کامل تمام تراکنش‌ها
- معرفی امکان «تأیید ساده پرداخت» (Simplified Payment Verification – SPV)
- گره: اعلام فهرستی از تراکنش‌های مرتبط با خود،
- به‌عنوان یک «گره سبک»
- امری که در تقابل با «گره‌های کامل» نگهدارنده کل داده‌های زنجیره

فیلتر بلوم

گره‌های سبک،

- فیلتر بلومی دارای فهرست تراکنش‌های مرتبط با آنها ساخته شده است
- ارسال به گره‌های کامل

گره کامل پیش از ارسال اطلاعات مربوط به یک تراکنش برای گره سبک،

بررسی فیلتر بلوم آن در ابتدا فیلتر بلوم

- کشف ارتباط یا بی‌ارتباطی

اگر «مثبت کاذب» نادیده گرفتن آن در گره سبک پس از دریافت پیام

فیلتر بلوم-ویژگی ها

امکان مثبت کاذب

عدم ذخیره مقادیر اصلی در فیلتر بلوم

اتکا به درهم‌سازهای محاسبه شده

ذخیره آنها در آرایه‌ای بیتی

با وجود فضا-کارآمدی چنین نمایشی

▪ منجر به موقعیت‌هایی تایید وجود برخی عناصر که غیر عضو

به دلیل تصادم‌های درهم

عدم وجود دانش قبلی در عملیات آزمون جهت مقایسه مقدار آزمون

فیلتر بلوم-ویژگی‌ها

اصل فیلتر بلوم: هنگام استفاده از فهرست یا مجموعه، و با ارزشی فضا، در صورت امکان تا زمانی که اثر مثبت کاذب چندان بالا نرفته یا آستانه‌ای را رد نکرده است از فیلتر بلوم بهره گرفته شود.

فیلتر بلوم-ویژگی ها

نادر بودن رخداد مثبت کاذب

امکان تخمین احتمال آن

مقدار کران پائین احتمال

$$\text{pr}_{fp} \approx \left(1 - e^{-\frac{kn}{m}}\right)^k$$

فیلتر بلوم-ویژگی ها

اگر t کسری از بیت‌های فیلتر بلوم پس از انجام n درج هنوز برابر صفر k تعداد توابع درهم‌سازی

احتمال رخداد مثبت کاذب (f) برابر با احتمال مشاهده k بیت یک

- به دلیل تایید حضور نیاز به برقراری k بیت با مقدار ۱

دستیابی به k مقدار ۱ می‌تواند نتیجه جست‌وجوی موفق برای عضو واقعاً درج شده

در صورت انتخاب تصادفی پرس‌وجوها از دامنه‌ای بسیار بزرگ‌تر از مجموعه داده

- آنگاه احتمال مثبت واقعی، کسری ناچیز از این مقدار
- ناشناختگی مقدار t پیش از آن که تمام عملیات درج
- وابستگی به خروجی درهم‌ساز

امکان کار با احتمال p

▪ احتمال صفر بودن یک بیت ثابت پس از انجام n درج

▪ مناسب جهت تقریب احتمالاتی درصد بیت‌های صفر موجود در فیلتر (t) $p = \left(1 - \frac{1}{m}\right)^{nk}$

فیلتر بلوم-ویژگی ها

فهم دلیل:

▪ آغاز از فیلتر بلوم خالی

▪ پس از آن که نخستین تابع درهم‌سازی $h1$ یکی از بیت‌ها را برابر ۱ قرار می‌دهد،

▪ احتمال این‌که یک بیت دلخواه در فیلتر بلوم مقدار ۱ داشته باشد برابر $1/m$ است و احتمال صفر بودن آن برابر $1 - \frac{1}{m}$

• پس از آن که همه k تابع درهم در درج نخست بیت‌های مربوطه را ۱ کردند، احتمال این‌که آن بیت مفروض همچنان

$$\left(1 - \frac{1}{m}\right)^k$$

پس از درج کل n مقدار، احتمال صفر بودن بیت مفروض

$$\left(1 - \frac{1}{m}\right)^{nk}$$

$$\left(1 + \frac{1}{x}\right)^x \approx e$$

$$p \approx e^{-\frac{nk}{m}}$$

میزان انحراف درصد واقعی صفرها انحراف از این مقدار؟

فیلتر بلوم-ویژگی ها

با استفاده از کران‌های چرنوف

▪ قضیه‌ای محدودساز احتمال انحراف شدید متغیر تصادفی از میانگینش

امکان نمایش کسر تمرکز مقدار کسر صفرها در فیلتر بلوم حول مقدار میانگین

کران چرنوف برای متغیرهای تصادفی X برابر جمع متغیرهای دنباله‌ای مستقل از مقادیر دوتائی

$$X = \sum_{i=1}^n X_i$$

X تعداد کل بیت‌های صفر موجود در فیلتر بلوم

▪ اگر $X_i = 0$ بیت i -ام مقدار ۱ و $X_i = 1$ در غیر این صورت

فیلتر بلوم-ویژگی ها

استفاده از کران چرنوف جهت نمایش عدم انحراف قابل توجه مقدار X از میانگینش

- عدم استقلال متغیرهای X_i در مسئله و وجود اندکی همبستگی منفی
- البته مطلوب تر

- در صورت برابر شدن بیتی برابر یک کاهش احتمال ۱ شدن سایر بیت ها

بیان کلی کران بالای چرنوف

- با μ به عنوان میانگین متغیر تصادفی X

$$Pr[X \geq (1 + \delta) \mu] \leq e^{-\mu \frac{\delta^2}{3}}$$

فیلتر بلوم-ویژگی ها

اعمال آن در مسئله

$$\mu = E[X] = mp = me^{-\frac{nk}{m}}$$

$$\delta = 1$$

احتمال این که X بیش از دو برابر مقدار میانگین خود انحراف پیدا کند

$$Pr[X \geq 2\mu] \leq e^{-\frac{\mu}{3}}$$

ادعا $nk = \theta(m)$

به معنای مقدار نمایی کوچک بودن احتمال انحراف X از میانگین به اندازه بیش از ضریب ۲،
 $e^{-\theta(m)}$.

بنابراین، p تقریب بسیار خوبی برای t یا درصد صفرهای فیلتر بلوم

فیلتر بلوم-ویژگی‌ها

با مشاهده کران، با تعداد مفروض و ثابت عناصر مورد انتظار n ،

وابستگی احتمال مثبت کاذب به انتخاب k و m

نیاز به بررسی و سبک-سنگینی واضح بین طول فیلتر، تعداد توابع درهم‌ساز، و احتمال چنین رویدادهایی

بدترین حالت هنگام پر بودن فیلتر

▪ به معنای یک شدن همه بیت‌ها

هر جستجو پاسخ مثبت (کاذب)

▪ انتخاب m وابسته به تعداد (تخمینی) موردانتظار n

▪ ضرورت بزرگ بودن m نسبت به n

نحوه محاسبه طول فیلتر m

▪ با توجه به احتمال مثبت کاذب pr_{fp} و تعداد مورد انتظار n

$$m = -\frac{n \ln pr_{fp}}{(\ln 2)^2}$$

فیلتر بلوم-ویژگی ها

نتیجه: در صورت اصرار بر حفظ میزان احتمال مثبت کاذب

- نیاز به رشد خطی فیلتر با نسبتی خطی از تعداد عناصر

امکان تنظیم احتمال مثبت کاذب با داشتن نسبت $\frac{m}{n}$

- به معنای تعداد بیت‌های تخصیص داده شده به ازای هر مقدار
- تنظیم با انتخاب تعداد توابع درهم‌ساز k

انتخاب بهینه k با هدف کمینه کردن احتمال مثبت کاذب

$$k = \frac{m}{n} \ln 2$$

تعداد بهینه توابع درهم‌ساز k حدود ۰,۷ تعداد بیت‌ها به ازای هر عنصر

ترجیح مقادیر زیربهینه کوچکتر به دلیل نیاز به عدد صحیح بودن k

فیلتر بلوم-ویژگی ها

برخی از راه‌حل‌های نزدیک به بهینه که با استفاده فراوان

pr_fp	m/n	k
0.0561	6	4
0.0215	8	6
0.00314	12	8
0.000458	16	11

فیلتر بلوم-ویژگی ها

مثال - تخمین پارامترها

احتمال مثبت کاذب $pr_{fp} = 1\%$

- نیاز است فیلتر ده برابر طولانی تر از تعداد عناصر مورد انتظار n
- استفاده از شش تابع درهم ساز
- عدم وابستگی طول فیلتر به اندازه خود مقادیر ورودی و ثابت ماندن برای مقادیر با ماهیت های مختلف

فیلتر بلوم-ویژگی ها

فیلترهای بلوم به مثابه تعمیمی از جدول های درهم

فیلتر با یک تابع درهم ساز معادل جدول درهم ساز

امکان ثابت نگه داشتن احتمال مثبت کاذب با استفاده از توابع درهم ساز متعدد در فیلترهای بلوم حتی برای تعداد ثابتی از بیت ها به ازای هر عنصر

▪ عدم امکان در جدول های درهم

فیلتر بلوم-ویژگی ها

ب- منفی کاذب امکان ناپذیر

▪ اگر خروجی فیلتر بلوم عدم عضویت عنصری باشد، عدم وجود قطعی مقدار مذکور

$$pr_{fn} = 0$$

فیلتر بلوم-ویژگی ها

ج- فیلتر بلوم مناسب برای استفاده در حافظه داخلی
کاهش احتمال مثبت کاذب با تخصیص حافظه بیشتر

تمایل به ایجاد فیلترهای بزرگتر (با m بزرگتر)

به محض بزرگ شدن بیش از حد، لاجرم و به اجبار انتقال به دیسک

چنین انتقالی با مشکل ناشی از طراحی

نیاز به دسترسی تصادفی با تولید شاخص‌های تصادفی در توابع درهم‌ساز با توزیع یکنواخت

ناکارآمدی نیاز مذکور برای دیسک‌های با صفحه‌های گردان و نوک‌های خواندن نوشتن متحرک

▪ همین‌طور برای دستگاه‌های ذخیره‌سازی حالت جامد

فیلتر بلوم-ویژگی ها

مثال - حافظه مورد نیاز: با توجه به معادله تخمین m ، جهت مدیریت یک میلیارد عنصر و حفظ احتمال رویدادهای مثبت کاذب حدود دو درصد و استفاده از تعداد بهینه‌ای از توابع درهم‌ساز،

فیلتر بلوم-ویژگی ها

مثال - حافظه مورد نیاز: با توجه به معادله تخمین m ، جهت مدیریت یک میلیارد عنصر و حفظ احتمال رویدادهای مثبت کاذب حدود دو درصد و استفاده از تعداد بهینه‌ای از توابع درهم‌ساز،

$$نیاز به انتخاب فیلتری به طول $m = -10^9 \times \frac{\ln 0.02}{(\ln 2)^2} \approx 8.14 \times 10^9$ بیت انتخاب کنیم$$

فیلتر بلوم-ویژگی ها

مثال - حافظه مورد نیاز: با توجه به معادله تخمین m ، جهت مدیریت یک میلیارد عنصر و حفظ احتمال رویدادهای مثبت کاذب حدود دو درصد و استفاده از تعداد بهینه‌ای از توابع درهم‌ساز،

$$m = -10^9 \times \frac{\ln 0.02}{(\ln 2)^2} \approx 8.14 \times 10^9$$

نیاز به انتخاب فیلتری به طول

تقریباً برابر با یک گیگابایت حافظه

فیلتر بلوم-ویژگی ها

ادغام دو فیلتر بلوم متفاوت با طول یکسان

- تنها در صورتی که دارای توابع درهم‌ساز یکسان
- تحویل ادغام به عملیات یای بیتی
- حاصل آن با فیلتر بلوم حاصل از اتحاد آن دو مجموعه از عناصر کاملاً متناظر

اشتراک دو فیلتر بلوم نیز ممکن است و با **and** بیتی صورت می‌پذیرد (چرا؟)

- با این حال، نتیجه امکان تصادم بیشتر

در صورت اتمام فضای فیلتر بلوم و نیاز به افزایش فضای آن

- ناممکن بودن بدون محاسبه مجدد تمام درهم‌هایی که قبلاً در فیلتر قرار گرفتند
- ناممکن در کاربردهای داده بزرگ

فیلتر بلوم-ویژگی ها

مثال - اشتراک کش

فهرستی از پروکسی های کش توزیع شده (P_1, P_2, \dots, P_n) در شبکه ای که کش های خود را به اشتراک می گذارند

در صورت ذخیره محتوای URL در خواستی روی پروکسی P_i

- برگرداندن آن در پروکسی بدون تماس واقعی با سرور راه دور
- در غیر این صورت، بازیابی محتوا، ذخیره محلی، و ارسال به مشتری

با هدف به کمینه سازی ترافیک شبکه و توزیع ذخیره سازی، می توان مسیریابی را در شبکه پروکسی تنظیم و سعی کنیم به سمت پروکسی ای هدایت کنیم که قبلاً محتوا را ذخیره کرده است. در غیر این صورت، با سرور راه دور تماس بگیریم. از آنجایی که درخواست های مشتری می توانند به هر پروکسی بیایند، مشکل اشتراک گذاری مسیریابی در هر پروکسی وجود دارد و زمانی که تغییر می کند، تبادل آن در شبکه یا ادغام، در صورت لزوم، وجود دارد.

فیلتر بلوم

- انتخابی طبیعی برای ذخیره چنین فهرست های مسیریابی و انجام پرسش های عضویت سریع
- به دلیل اندازه کوچک خود به راحتی در شبکه قابل انتقال

فیلتر بلوم-ویژگی ها

رویداد مثبت کاذب، در این مورد، این است که فرضی یکی از پروکسی‌ها مثلاً P_i فرض کند که پروکسی دیگری P_j ممکن است محتوای URL درخواستی را داشته باشد، اما در واقع ندارد. P_i ترافیک را به P_j هدایت می‌کند و از آن می‌خواهد محتوا را برگرداند، بنابراین P_j باید با سرور راه دور تماس بگیرد. در نتیجه، مقداری ترافیک شبکه اضافی تولید می‌کند و کپی‌های محلی اضافی را برای چنین محتوایی ذخیره می‌کند، که کاملاً قابل قبول است.

فیلتر بلوم-ویژگی‌ها

د- حذف امکان‌پذیر نیست.

برای حذف عنصری خاص از فیلتر بلوم نیاز به پاکسازی بیت‌های k متناظر آن در آرایه بیتی

امکان مطابقت یک بیت واحد با چند مقدار

- به دلیل تصادم‌های درهم و بیت‌های مشترک بین مقادیر

معرفی چند بهبود جهت پشتیبانی حذف

- لنگیدن یک پای ماجرا

- تحمیل هزینه‌هایی از نظر فضا و سرعت

به همین دلیل است که فیلتر بلوم کلاسیک بسیار سریع و کارآمد از نظر فضا است. نکته مثبت این است که عدم پشتیبانی از حذف در بسیاری از برنامه‌های کاربردی دنیای واقعی مشکلی ایجاد نمی‌کند، اما اگر واقعاً به آن نیاز باشد، چاره‌ای جز به بهبود فیلتر بلوم مانند «فیلتر بلوم شمارنده» نیست.

شمارش عناصر منحصر به فرد در فیلتر

روشی برای تخمین تعداد عناصر منحصر به فرد نمایه شده در فیلتر

اس. جاشوا سامیاس و پیر بالدی

▪ در واقع، بهبودی بر «الگوریتم شمارش خطی» است که در بخش‌های بعد بحث

با استفاده از اطلاعات مربوط به تعداد بیت‌های تنظیم شده در فیلتر و احتمال تنظیم هر بیت

▪ عرضه معادله ساده‌ای برای تقریب تعداد عناصر موجود در فیلتر

▪ عدم تغییر بیت‌های تنظیم شده در ازای افزودن دو مقدار برابر به فیلتر

▪ ارائه چنین تقریبی تخمینی برای تعداد عناصر منحصر به فرد (که به عنوان تعداد اصلی شناخته می‌شود)

شمارش عناصر منحصر به فرد در فیلتر

الگوریتم - شمارش عناصر منحصر به فرد در فیلتر بلوم

ورودی: فیلتر بلوم به طول m با k تابع درهم ساز

خروجی: تعداد عناصر منحصر به فرد در فیلتر

N برابر شمارش Bloomfilter[j] برای $j = 1, \dots, m$

اگر $N < k$

برگرداندن صفر

اگر $N == k$

برگرداندن یک

اگر $N == m$

برگرداندن m/k

برگرداندن $-\frac{m}{k} \cdot \ln\left(1 - \frac{N}{m}\right)$

شمارش عناصر منحصر به فرد در فیلتر

احتمال صفر بودن بیتی پس از درج n عضو جدید برابر است با

$$\text{pr}(\text{bit} = 0) = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-kn/m}$$

در نتیجه، میانگین نسبت یک‌ها برابر است با

$$\frac{E[N]}{m} = 1 - e^{-kn/m}$$

و N تعداد بیت‌های متناظر یک در فیلتر

جهت تقریب n از N می‌توان از $N/m = 1 - e^{-kn/m}$ مقدار زیر را بدست آورد:

$$e^{-kn/m} = 1 - \frac{N}{m}$$

شمارش عناصر منحصر به فرد در فیلتر

$$e^{-kn/m} = 1 - \frac{N}{m}$$
$$-\frac{kn}{m} = \ln\left(1 - \frac{N}{m}\right)$$
$$n = -\frac{m}{k} \ln\left(1 - \frac{N}{m}\right)$$

که دقیقا برابر مقداری است که الگوریتم وقتی $0 < N < m$ برمی گرداند. حالت های خاص آن عبارت است از الف- $N < k$ که غیرممکن است مگر $n = 0$. ب- $N = k$ که نشان می دهد احتمالا یک عضو افزوده شدن. ج- $N = m$ یا تمامی بیت ها برابر ۱ است. در نتیجه ظرفیت بیشینه یا m/k عضو دارد.

معادله $-\frac{m}{k} \ln(1 - N/m)$ تقریبی مناسب از n در مواقعی است که فیلتر اشباع نشده است.

فیلتر بلوم شمارنده

معروف‌ترین بهبود بر فیلتر بلوم کلاسیک

- پشتیبانی از حذف
- لی فن، پی کائو، جوسارا آلمیدا، و آندری برودر در سال ۱۳۷۹
- با تکیه بر الگوریتم فیلتر بلوم کلاسیک،

آرایه‌ای از m شمارنده $\{C_j\}_{j=1}^m$

- متناظر با هر بیت در آرایه فیلتر

با افزوده شدن هر عنصر افزایش یک واحدی مدخل‌های متناظر
در نتیجه، فراهم کردن تقریبی از تعداد دفعاتی مشاهده هر عنصر در فیلتر

داده‌ساختار مرتبط

- شامل آرایه‌ای بیتی و آرایه‌ای از شمارنده‌ها هر دو به طول m
- همگی مقداردهی اولیه به صفر

فیلتر بلوم شمارنده

هنگام درج مقداری جدید

- محاسبه موقعیت‌های بیت متناظر آن
- افزایش شمارنده متناظر برای هر موقعیت
- تغییر آرایهٔ بیتی از صفر به یک با تغییر شمارنده از صفر به یک

فیلتر بلوم شمارنده

الگوریتم - افزودن عنصر به فیلتر بلوم شمارنده

ورودی: عنصر $x \in D$

ورودی: فیلتر بلوم شمارنده با m شمارنده $\{C_j\}_{j=1}^m$ و k تابع درهم‌ساز $\{h_i\}_{i=1}^k$

برای i از ۱ تا k انجام بده

$$j \leftarrow h_i(x)$$

$$C_j \leftarrow C_j + 1$$

اگر $C_j == 1$

$CountingBloomFilter[j] \leftarrow 1$

فیلتر بلوم شمارنده

عملیات آزمون دقیقاً همانند فیلتر بلوم
▪ عدم نیاز به بررسی شمارندهها

به دلیل آرایه‌های بیتی یکسان فیلترها
▪ دارای زمان برابر برای آزمون عنصری با الگوریتم کلاسیک

فیلتر بلوم شمارنده

الگوریتم- آزمایش عنصر در فیلتر بلوم شمارنده

ورودی: عنصر $x \in D$

ورودی: فیلتر بلوم شمارنده با k تابع درهم‌ساز $\{h_i\}_{i=1}^k$

خروجی: غلط در صورت نیافتن عنصر و درست در صورت وجود عنصر

برای i از ۱ تا k انجام بده

$$j \leftarrow h_i(x)$$

اگر $CountingBloomFilter[j] \neq 1$

برگرداندن غلط

برگرداندن درست

فیلتر بلوم شمارنده

حذف برعکس درج

حذف عنصر x ،

- محاسبه تمامی k مقدار درهم $\{h_i\}_{i=1}^k$
- کاهش شمارنده‌های متناظر

در صورت تغییر مقدار شمارنده از یک به صفر

- حذف بیت متناظر در آرایه بیتی

فیلتر بلوم شمارنده

الگوریتم - حذف عنصر از فیلتر بلوم شمارنده

ورودی: عنصر $x \in D$

ورودی: فیلتر بلوم شمارنده با m شمارنده $\{C_j\}_{j=1}^m$ و k تابع درهم‌ساز $\{h_i\}_{i=1}^k$

برای i از ۱ تا k انجام بده

$$j \leftarrow h_i(x)$$

$$C_j \leftarrow C_j - 1$$

$$C_j == 0 \text{ اگر}$$

$$\text{CountingBloomFilter}[j] \leftarrow 0$$

فیلتر بلوم شمارنده

فرض الگوریتم حذف عنصر

- وجود مقدار x در فیلتر
- یا امکان وجود

نیاز به الگوریتم آزمون عضویت پیش از کاهش شمارنده‌های متناظر

وارث برحق!

- به ارث بردن تمام ویژگی‌های فیلتر بلوم معمول از جمله تخمین خطای مثبت کاذب و معادلات انتخاب بهینه m و k

فیلترهای بلوم شمارنده بسیار بزرگتر از فیلترهای بلوم کلاسیک

- به دلیل نیاز به حافظه اضافی برای شمارنده‌ها ولو اینکه بیشتر مقادیر برابر صفر
- اهمیت تخمین میزان بزرگ شدن چنین شمارنده‌هایی
- وابستگی اندازه آنها به طول فیلتر m و تعداد توابع درهم‌ساز k

با فرض هر شمارنده C دارای سطح ظرفیت N

- احتمال فراتر رفتن مقدار از N (به معنای سرریزی) با طول m با انتخاب بهینه k برابر مقدار

$$pr(\max(C) \geq N) \leq m \times \left(\frac{e \ln 2}{N}\right)^N$$

فیلتر بلوم شمارنده

مثال - شمارنده‌های ۴ بیتی $N = 16$ ، احتمال سرریز با معادله بالا برابر است با

فیلتر بلوم شمارنده

مثال - شمارنده‌های ۴ بیتی $N = 16$ ، احتمال سرریز با معادله بالا برابر است با

$$pr(\max(C) \geq 16) \leq m \times 1.37 \times 10^{-15}$$

فیلتر بلوم شمارنده

مثال - شمارنده‌های ۴ بیتی $N = 16$ ، احتمال سرریز با معادله بالا برابر است با

$$pr(\max(C) \geq 16) \leq m \times 1.37 \times 10^{-15}$$

به عبارت دیگر، اگر ۴ بیت به ازای هر شمارنده تخصیص دهیم، احتمال سرریز برای مقادیر عملی m (به عنوان مثال، چند میلیارد موقعیت بیتی) در طول درج اولیه در فیلتر بسیار کم است.

فیلتر بلوم شمارنده

مثال - شمارنده‌های ۴ بیتی $N = 16$ ، احتمال سرریز با معادله بالا برابر است با

$$pr(\max(C) \geq 16) \leq m \times 1.37 \times 10^{-15}$$

به عبارت دیگر، اگر ۴ بیت به ازای هر شمارنده تخصیص دهیم، احتمال سرریز برای مقادیر عملی m (به عنوان مثال، چند میلیارد موقعیت بیتی) در طول درج اولیه در فیلتر بسیار کم است.

پس از حذف و درج‌های متعدد، احتمال می‌تواند کمی بیشتر شود، اما هنوز برای استفاده عملی به اندازه کافی کم است.

فیلتر بلوم شمارنده

جلوگیری از سرریز حسابی

- یا افزایش شمارنده‌ای که قبلاً حداکثر مقدار ممکن را دارد

هر شمارنده در آرایه به اندازه کافی بزرگ باشد تا ویژگی‌های فیلترهای بلوم حفظ شود.

- در عمل، شمارنده از ۴ یا چند بیت تشکیل شده است و فیلتر بلوم شمارنده

- نیاز به چهار برابر فضای بیشتر نسبت به فیلتر بلوم کلاسیک

فرارفتن مقدار شمارنده چهار بیتی از مقدار ۱۵

- تحویل آن «به حد بالا» و باقی ماندن مقدار در ۱۵

امکان تولید پاسخ منفی پس از حذف‌های زیاد

- صفر شدن شمارنده زمانی که نباید صفر شود

- اما احتمال چنین زنجیره‌ای از رویدادها آنقدر کم است که راه‌اندازی مجدد برنامه و ایجاد دوباره فیلتر بسیار محتمل‌تر

فیلتر بلوم شمارنده

با این حال، امکان طراحی نسخه پیچیده‌تری از فیلتر بلوم شمارنده با شمارنده‌های کوچکتر

- مثلاً دو بیتی

- در عوض با اتخاذ رویکردی مشابه آدرس‌دهی بسته در جدول‌های درهم و معرفی جدول درهم ثانویه برای مدیریت شمارنده‌های سرریز

بنابراین، فیلتر بلوم شمارنده به دلیل احتمال مثبت خطا به دلیل سرریز شدن شمارنده صرفاً از حذف درست به صورت احتمالی پشتیبانی می‌کند.

با وجود ایرادهای مورد اشاره، فیلترهای بلوم شمارنده در **Apache Hadoop** و **Apache Spark** در برنامه‌های **MapReduce** برای تسریع پردازش مجموعه‌های داده عظیم روی خوشه‌های بزرگ با کمک به کاهش حجم سخت مورد استفاده‌اند.

فیلتر خارج قسمت

مایکل بندر و همکاران در سال ۱۳۹۰ معرفی فیلتر خارج قسمت

- پشتیبانی از عملیات اصلی فیلترهای بلوم
- بهبود محلی سازی داده
- نیاز به تعداد کمی از دسترسی های دیسک پیوسته

دارای عملکرد فضا و زمان مشابه

پشتیبانی از حذف و امکان تغییر اندازه یا ادغام پویا

نام این داده ساختار ریشه از خارج قسمت نتیجه عملیات تقسیم

فیلتر خارج قسمت

فیلتر خارج قسمت در ساده ترین شکل خود

- جدول درهم با استفاده از کاوش خطی

به جای ذخیره سازی مقدارهای کلید در خانه ها، مانند جدول درهم خطی کلاسیک،

- ذخیره مقدارهای درهم در فیلتر خارج قسمت

- اصطلاح اثرانگشت به جای درهم

- ذخیره تکه ای از هر درهم در فیلتر خارج قسمت

- با امکان بازسازی کامل مقدار درهم به طور قابل اعتماد

فیلتر خارج قسمت

در «طیف دقت»

▪ فیلتر خارج قسمت جایی بین جدول درهم و فیلتر بلوم

درهم شدن دو مقدار کلید متمایز به یک اثرانگشت یکسان

▪ فیلتر خارج قسمت قادر به تشخیص تفاوت بین آن‌ها به شیوه‌ای که جدول درهم‌سازی می‌تواند، نخواهد بود.

▪ اما اگر دو مقدار کلید به اثرانگشت‌های متفاوتی درهم شوند، فیلتر خارج قسمت قادر خواهد بود آن‌ها را از یکدیگر تشخیص دهد.

▪ این وضعیت در مورد فیلترهای بلوم صادق نیست،

▪ جایی که پرس و جو بر روی یک مقدار کلید با مجموعه‌ای منحصر به فرد از k درهم ممکن است یک مثبت کاذب ایجاد نماید.

▪ فیلتر خارج قسمت کارکردی مشابه فیلتر بلوم با طراحی بسیار متفاوت

▪ با استفاده از اثرانگشت‌های طولی‌تر، فیلترهای خارج قسمت می‌توانند میزان مثبت کاذب را کاهش دهند

▪ اما اثرانگشت‌های طولی‌تر ممکن است فضای زیادی نیز مصرف کنند.

فیلتر خارج قسمت

صناعات و ترفندهای مختلفی مورد استفاده فیلتر خارج قسمت برای ذخیره سازی فشرده اثرانگشتها قابلیت بازسازی اثرانگشت کامل زمانی مفید واقع می شود که بخواهیم مقادارها را حذف کنیم.

- فیلترهای خارج قسمت می توانند به طور کارآمد حذف را انجام دهند.

فیلتر خارج قسمت قادر به تغییر اندازه

- ادغام دو فیلتر خارج قسمت در یک فیلتر خارج قسمت بزرگ تر

ادغام کارآمد، تغییر اندازه کارآمد، و حذف کارآمد

مفید بودن این ویژگی ها به ویژه در سیستم های توزیع شده پویا

درک و پیاده سازی فیلترهای خارج قسمت نسبت به فیلترهای بلوم اندکی پیچیده تر

- از طرف دیگر، به دلیل نیاز به دسترسی تصادفی چندباره به هر عملیات، کارایی فیلتر بلوم در صورت عدم گنجایش در حافظه اصلی و نیاز به حافظه جانبی سخت افت می کند و نامناسب می شود.

فیلتر خارج قسمت

طراحی فیلتر خارج قسمت

- توضیح خارج قسمت گیری
- توصیف چگونگی فیلتر خارج قسمت از بیت های متاداده همراه با خارج قسمت گیری برای صرفه جویی در فضا

از راه های به نظر آوردن فیلتر خارج قسمت

- ، مشاهده آن به عنوان بازی برای صرفه جویی یک بیت اینجا و آنجا
- به کارگیری ترفندهای مختلف در این راستا

فیلتر خارج قسمت تنها ساختمان داده از این نوع نیست،

- اما برخی از ترفندهایی که در آن استفاده می شود می توانند به طور کلی در درک ساختمان های داده مشابه مبتنی بر جدول های درهم فشرده مفید واقع شوند.

خارج قسمت گیری

خارج قسمت گیری

- مقدار درهم ورودی
- تقسیم به خارج قسمت و باقیمانده

استفاده از خارج قسمت برای نمایه گذاری در سطل متناظر جدول درهم در فیلتر خارج قسمت

ذخیره باقیمانده در خانه متناظر

برای مثال با فرض درهم به طول h بیت و اندازه جدول m معادل 2 به توان q

- خارج قسمت مقداری که با q بیت چپ درهم تعیین می شود و باقیمانده نشان دهنده r بیت باقیمانده به تعداد h منهای q بیت

مثال - با فرض $m=32$ و $h=11$ باشد پس $q=5$ و $r=6$

- صرفه جویی به دلیل خارج قسمت گیری به ازای هر خانه معادل $q = 5$ بیت در مثال قبلی

فیلتر خارج قسمت

جهت بهبود محلی سازی فضایی،

- داده ساختار فیلتر خارج قسمت
- نمایش با جدول درهم ساز فشرده با آدرس دهی باز با $m = 2^q$ سطل
- ذخیره باقیمانده f_r در سطل نمایه گذاری شده با خارج قسمت f_q
- حل تصادم های احتمالی با کاوش خطی

فیلتر خارج قسمت

الگوریتم - تکنیک خارج قسمت گیری

ورودی: اثر انگشت f

خروجی: خارج قسمت f_q و باقیمانده f_r

$$f_r \leftarrow f \% 2^r$$

$$f_q \leftarrow \left\lfloor \frac{f}{2^r} \right\rfloor$$

برگرداندن f_r و f_q

فیلتر خارج قسمت

درج در فیلتر خارج قسمت چگونه

$h = \text{len}(f)$ ← تعداد بیت‌های در درهم یا اثر انگشت

$q = \log_2(m)$ ← فرض بر اندازه درهم توانی از دو

$r = h - q$

$f_q = \text{math.floor}(f / 2^{**} r)$ ← q بیت سمت چپ اثر انگشت

$f_r = f \% 2^{**} r$ ← r بیت سمت راست اثر انگشت

$\text{filter}(f_q) = \text{rem}$ ← ذخیره باقی مانده در سطل متناظر خارج قسمت

فیلتر خارج قسمت

اگر هیچ دو اثرانگشتی دارای خارج قسمت یکسان نباشند

- به معنای عدم تصادم
- در این صورت، اشغال هر باقیمانده در سطل مخصوص به خود
- بازسازی درهم کامل از طریق الحاق بازنمایی دودویی شماره سطل b با بازنمایی دودویی باقیمانده ذخیره شده در سطل b

با توجه به باقیمانده f_r در سطل f_q ، بازسازی منحصر به فرد اثر انگشت کامل

$$f = f_q \times 2^r + f_r$$

کمک این ویژگی به تغییر اندازه و ادغام

- عدم امکان بازسازی مقادیر اصلی
- یادآوری جایی بین جدول‌های درهم نگهدارنده مقادیرهای کلید واقعی و فیلترهای بلوم بدون امکان بازسازی مقدارمتناظر درهم‌ها

فیلتر خارج قسمت

با این حال

- رایج بودن تصادم در جدول‌های درهم
- حل تصادم در فیلترهای خارج قسمت با استفاده از گونه‌ای از کاوش خطی

امکان رخداد پیچیدگی‌هایی

- در نتیجه انتقال برخی از باقیمانده‌ها از سطل اصلی خود به سمت خانه‌های بعدی بر اثر کاوش خطی
- از دست رفتن ارتباط بین خارج قسمت و باقیمانده

برای بازسازی درهم کامل،

- نیاز فیلتر خارج قسمت به ازای هر خانه از سه بیت متاداده
- سه بیت در مقابل صرفه‌جویی حدود ۲۰ تا ۳۰ بیت به ازای هر خانه بابت خارج قسمت در جدول‌های درهم بزرگ‌تر
- بهائی اندکی

فیلتر خارج قسمت

چند اصطلاح

مفهوم اجرا یا Run:

- هر اجرا توالی پیوسته‌ای از خانه‌های فیلتر خارج قسمت
- ذخیره باقی مانده‌ها (یا اثرانگشت‌ها) با خارج قسمت یکسان
- ذخیره همه اثرانگشت‌های دارای تصادمی روی یک سطل مشخص، به صورت پشت سرهم و به ترتیب صعودی باقی مانده‌ها

به علت وقوع تصادم‌ها و عمل هل دادن یا انتقال هنگام درج
▪ امکان آغاز اجرا از فاصله‌ای بسیار دورتر از مدخل سطل متناظر خود

مفهوم خوشه

- خوشه رشته‌ای پیوسته از یک یا چند اجرا
- رشته شامل خانه‌ها اشغال شده با باقیمانده
- قرارگیری اولین باقیمانده در خانه اصلی درهم شده خود
- به نام
- به نام مدخل لنگر و تکیه
- پایان خوشه یا با مدخلی خالی، یا با شروع خوشه‌ای دیگر

فیلتر خارج قسمت

انجام عملیات جست و جو در فیلتر باقیمانده و بازیابی اثرانگشت کامل

- نیاز به تفسیر باقیمانده‌ها با بیت‌های متاداده
- آغاز تفسیر از آغاز خوشه (موقعیت لنگر)
- نقش سه بیت متاداده موجود در هر خانه
- ۱. `bucket_occupied`: درهم‌شدن کلیدی روی سطل متناظر
 - مقدار ۱: درهم‌شدگی حداقل یک کلید به این مدخل
 - مقدار ۰: عدم درهم‌شدگی کلیدی به این سطل
- سخن کوتاه، بیت روشن‌گر تعیین خارج‌قسمت‌های حاضر در خوشه
- ۲. `run_continued`: باقیمانده‌های فعلی ادامه اجرای قبلی یا آغاز اجرای جدید
 - مقدار ۰: باقیمانده فعلی اولین عضو اجرا
 - مقدار ۱: تعلق باقیمانده فعلی به اجرای قبلی
- سخن کوتاه، مرزبندی دقیق اجراها با این بیت
- ۳. `is_shifted`: تشخیص و تبیین قرار گرفتن باقیمانده در خانه اصلی خود یا دچار جابجا شدگی به دلیل تصادم
 - مقدار ۰: درهم‌شدن باقیمانده در خانه اصلی
 - مقدار ۱: انتقال باقیمانده به سمت **خانه‌های بعدی**
- یافتن آغاز خوشه با این بیت

فیلتر خارج قسمت

- هر سطل شامل سه بیت متادادهٔ `bucket_occupied`، `run_continued`، و `is_shifted`
- همگی در ابتدا بدون مقدار
- نقشی مهم در مدیریت داده ساختار

فیلتر خارج قسمت

شکل - سطل در فیلتر خارج قسمت

f_q سطل

Bucket_occupied اشغال سطل	Run_continued ادامه اجرا	ls_shifted منتقل شده یا سر جای خود
f_r		

درج در فیلتر خارج قسمت

مثال جهت فهم بهتر

بررسی گام به گام فرایند درج مقدارهای v ، w و x

آغاز با فیلتر خارج قسمت در ابتدا خالی با ۳۲ خانه برابر با ۲ به توان ۵

- تنظیم هر خانه و همچنین هر سه بیت متاداده ابتدا به مقدار ۰

درج در فیلتر خارج قسمت

خارج قسمت / نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱				
۱۰۰۱۰				
۱۰۰۱۱				
۱۰۱۰۰				
۱۰۱۰۱				
۱۰۱۱۰				
۱۰۱۱۱				

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

۱ درج ۷:

- مقدار درهم $h(v)$ برابر با 10001001010
- سطلی که با خارج قسمت 10001 تعیین می شود قبلاً اشغال نشده است.
- بیت `bucket_occupied` را روی مقدار ۱ تنظیم و ذخیره باقیمانده را در خانه‌ای که با خارج قسمت آن تعیین می شود
- نیازی به هیچ اقدام اضافی روی سایر بیت‌ها نیست
- زیرا این مورد در حال حاضر هم آغاز یک اجرا و هم آغاز یک خوشه محسوب می شود

درج در فیلتر خارج قسمت

خارج قسمت / نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰۰۱۰				
۱۰۰۱۱				
۱۰۱۰۰				
۱۰۱۰۱				
۱۰۱۱۰				
۱۰۱۱۱				

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

۲ درج w :

مقدار درهم $h(w)$ برابر با ۱۰۰۱۱۱۰۰۱۱۱

مجدداً بیت `bucket_occupied` را در خارج قسمت ۱۰۰۱۱ روی مقدار ۱ تنظیم

ذخیره باقیمانده در خانه متناظر در دسترس

عدم تغییر بیت‌های دیگر

درج در فیلتر خارج قسمت

خارج قسمت/نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰۰۱۰				
۱۰۰۱۱	۱۰۰۱۱۱	۱	۰	۰
۱۰۱۰۰				
۱۰۱۰۱				
۱۰۱۱۰				
۱۰۱۱۱				

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

۳ درج x :

- مقدار درهم $h(x)$ برابر با 10011111111
- سطل متناظر با خارج قسمت 10011 هنگامی که سعی می‌کنیم بیت را روی مقدار 1 تنظیم کنیم از پیش اشغال شده است.
- بنابراین هر جا که باقیمانده را ذخیره کنیم باید بیت `run_continued` آن خانه را روی مقدار 1 تنظیم کنیم.
- خانه در سطل درهم‌شده اشغال می‌باشد.
- از این رو هر جا که باقیمانده را ذخیره کنیم باید بیت `is_shifted` آن خانه را نیز روی مقدار 1 تنظیم نماییم.
- با توجه به این که در ابتدای خوشه‌ای قرار داریم که فقط یک اجرا دارد (که خارج قسمت آن با خارج قسمت x برابر است)،
 - پیمایش به سمت پایین تا یافتن نخستین خانه در دسترس درون اجرا در سطل 10100
- ذخیره باقیمانده
- تنظیم مقدار 1 بیت‌های `run_continued` و `is_shifted`

درج در فیلتر خارج قسمت

خارج قسمت / نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰۰۱۰				
۱۰۰۱۱	۱۰۰۱۱۱	۱	۰	۰
۱۰۱۰۰	۱۱۱۱۱۱	۰	۱	۱
۱۰۱۰۱				
۱۰۱۱۰				
۱۰۱۱۱				

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

در حال حاضر فیلتر خارج قسمت

- دارای دو خوشه است که هر خوشه دارای یک اجرا
- در ادامه، درج چند مقدار دیگر

درج در فیلتر خارج قسمت

۴ درج y :

- مقدار درهم $h(y)$ برابر با 10100101101
- بیت `bucket_occupied` در خارج قسمت 10100 مقدار صفر داشت
- تغییر آن به مقدار 1
- بیت `run_continued` در خانه باقیمانده نهایی روی مقدار 0 تنظیم خواهد شد
- خانه در سطل درهم شده اشغال
- بیت `is_shifted` در خانه باقیمانده نهایی روی مقدار 1 تنظیم خواهد شد
- با شروع از ابتدای خوشه به دنبال نخستین مکان برای ذخیره اجرای جدید خود می گردیم.
- نخستین خانه در دسترس در خارج قسمت 10101 قرار دارد.
- از این رو باقیمانده را ذخیره کرده و بیت های متاداده را متناسب با آن تنظیم می کنیم.

درج در فیلتر خارج قسمت

خارج قسمت / نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰۰۱۰				
۱۰۰۱۱	۱۰۰۱۱۱	۱	۰	۰
۱۰۱۰۰	۱۱۱۱۱۱	۱	۱	۱
۱۰۱۰۱	۱۰۱۱۰۱	۰	۰	۱
۱۰۱۱۰				
۱۰۱۱۱				

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

۵ درج z :

- مقدار درهم $h(z)$ برابر با 10100111110
- بیت `bucket_occupied` سطل مربوط به z از پیش روی مقدار ۱ تنظیم شده است
- بیت `run_continued` در خانه نهایی روی مقدار ۱ خواهد بود
- خانه اصلی اشغال است
- بیت `is_shifted` در خانه نهایی روی مقدار ۱ خواهد بود
- با شروع از ابتدای خوشه محل اجرای مناسب خود را پیدا می‌کنیم.
- در امتداد اجرا به سمت پایین تا خارج قسمت 10110 پویش می‌کنیم تا باقیمانده z را ذخیره کنیم و بیت‌ها را متناسب با آن تنظیم می‌نماییم.

درج در فیلتر خارج قسمت

خارج قسمت / نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰۰۱۰				
۱۰۰۱۱	۱۰۰۱۱۱	۱	۰	۰
۱۰۱۰۰	۱۱۱۱۱۱	۱	۱	۱
۱۰۱۰۱	۱۰۱۱۰۱	۰	۰	۱
۱۰۱۱۰	۱۱۱۱۱۰	۰	۱	۱
۱۰۱۱۱				

اجرا
خوشه
اجرا

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

به دلیل ترتیب صعودی مقادیر درج شده، عمل درج نسبتاً ساده
امکان ترتیب مرتب شده درج‌ها در فیلتر خارج قسمت پیرنگی
▪ اما غالب نیست

همچنین، باید توانایی مدیریت پیرنگ‌هایی ممکن باشد که در آن اثرانگشت‌های درج شده به ترتیب دلخواه و
نامرتب می‌رسند.

درج در فیلتر خارج قسمت

درج خارج از ترتیب مرتب شده اثرانگشت‌ها

پس از یافتن اجرای صحیح برای باقیمانده‌ای که قرار است در آن درج شود

امکان جابجایی چندین مورد از آن اجرا و سایر مقدارهای موجود در آن خوشه به خانه‌های بعدی

نمونه درج یک مقدار با نام a با مقدار درهم $h(a)$ برابر با 10100000000

این مقدار متعلق به دومین اجرای دومین خوشه

- اشغال خانه‌های تعیین شده با خارج قسمت‌های 10101 و 10110

- جابجایی یک خانه کل اجرا با مقدار a به سمت پایین

- ذخیره در خانه 10101

- زیرا باقیمانده آن نخستین مقدار در ترتیب صعودی مرتب شده در آن اجرا محسوب می‌شود

- در نتیجه همچنین باعث ایجاد تغییرات در بیت‌های متاداده می‌گردد.

- باقیمانده‌ها باید در درون یک اجرا مرتب شده باشند و همچنین اجراها در میان خوشه‌ها مرتب باشند

- چرا؟

درج در فیلتر خارج قسمت

خارج قسمت/نشانی	باقیمانده	متاداده		
۱۰۰۰۰				
۱۰۰۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰۰۱۰				
۱۰۰۱۱	۱۰۰۱۱۱	۱	۰	۰
۱۰۱۰۰	۱۱۱۱۱۱	۱	۱	۱
۱۰۱۰۱	۰۰۰۰۰۰	۰	۰	۱
۱۰۱۱۰	۱۰۱۱۰۱	۰	۱	۱
۱۰۱۱۱	۱۱۱۱۱۰	۰	۱	۱

اجرا
خوشه
اجرا

سطل اشغال شده: مقداری به این سطل درهم شده است

اجرا ادامه دار: مقدار فعلی ادامه اجرایی است

منتقل شده: این مدخل انتقالی است

درج در فیلتر خارج قسمت

اهمیت کوچک بودن اندازه خوشه‌ها

- نیاز بالقوه به جابجایی کل یک خوشه از مقادیر در حین درج یا حذف
- همچنین رمزگشایی کل یک خوشه در حین انجام عملیات جستجو

هر چه فضای خالی بیشتری

- کاهش احتمال تشکیل خوشه‌های بزرگ که عملیات درج و جستجو نیاز به پویس و رمزگشایی در آن‌ها داشته باشد
- درست مانند جدول‌های درهم معمولی با کاوش خطی،
- فیلترهای خارج قسمت زمانی سریع‌تر کار می‌کنند که ضریب بار در محدوده ۷۵ تا ۹۰ درصد نگه داشته شود

فیلتر خارج قسمت

فیلتر خارج قسمت مجموعه داده $D = \{x_1, x_2, \dots, x_n\}$ با ذخیره اثر انگشتی p -بیتی برای هر یک از اعضاء مجموعه داده وابستگی به تابع درهم‌ساز برای تولید چنین اثر انگشت‌هایی

جهت پشتیبانی از تصادفی‌سازی مناسب

- نیاز به تابع درهم‌ساز با تولید اثر انگشت‌هایی با توزیع یکنواخت و مستقل
- هر اثر انگشت f در الگوریتم تقسیم به دو بخش
 - استفاده از تکنیک خارج‌قسمت‌گیری
 - q بیت با بیشترین اهمیت (خارج‌قسمت)
 - r بیت با کمترین اهمیت (باقیمانده)
- پیشنهاد دونالد کنوٹ

فیلتر خارج قسمت

بهبود محلی سازی فضایی

- داده ساختار فیلتر خارج قسمت
- نمایش با جدول درهم ساز فشرده با آدرس دهی باز با $m = 2^q$ سطل
- ذخیره باقیمانده f_r در سطل نمایه شده با خارج قسمت f_q
- حل تصادم های احتمالی با کاوش خطی

وجود باقیمانده f_r در سطل f_q ،
بازسازی منحصر به فرد اثر انگشت کامل

$$f = f_q \times 2^r + f_r$$

دو اثر انگشت متفاوت f و f' با خارج قسمت یکسان ($f_q = f'_q$)

- تصادم نرم
- حل با کاوش خطی
- پیاده سازی با ذخیره پشت سر هم تمام باقیمانده های اثر انگشت ها با خارج قسمت یکسان در اجرا
- امکان انتقال باقیمانده از مکان اصلی خود به خانه های بعدی و ذخیره در سطل بعدی
- حلقه در انتهای آرایه

فیلتر خارج

	۱	ورودی: نمایه سطل i
	۲	ورودی: فیلتر خارج قسمت به طول m
	۳	$qbl \leftarrow QF[i]$
	۴	$k \leftarrow i + 1$
	۵	تازمانی که درست (True) انجام بده
	۶	اگر $k \geq m$ آنگاه
	۷	$k \leftarrow \cdot$
	۸	اگر $QF[k] = \phi$ آنگاه
	۹	$QF[k] \leftarrow qbl$
هر مقدار منتقل شده ادامه اجراست	۱۰	$QF[k].run_continued \leftarrow 1$
در موقعیت اصلی خود نیست	۱۱	$QF[i] \leftarrow \phi$
	۱۲	$QF[i].bucket_occupied \leftarrow \cdot$
	۱۳	برگرداندن QF
	۱۴	ورنه
	۱۵	$knoni \leftarrow QF[k]$
	۱۶	$QF[k] \leftarrow qbl$
تبدیل به ادامه اجرا	۱۷	$QF[k].run_continued \leftarrow 1$
منتقل شده	۱۸	$QF[k].is_shifted \leftarrow 1$
	۱۹	$qbl \leftarrow knoni$
	۲۰	$k \leftarrow k + 1$

فیلتر خارج قسمت

خوشه

- دنباله‌ای از یک یا چند اجرای متوالی بدون سطل خالی در بین آنها

قرارگیری همه خوشه‌ها بلافاصله پیش از سطل خالی

- عدم مقداردهی بیت `is_shifted` اولین مقدار آن

فیلتر خارج قسمت

ذخیره جدول درهم داخلی به صورت فشرده در آرایه‌ای

- کاهش حافظه موردنیاز کاهش
- بهبود محلی‌سازی داده

منجر به پیچیده‌تر شدن مدیریت

تابع اسکن یا پیمایش برای یافتن اجرا

- گام‌برداری عقب‌رو از سطل متعارف f برای یافتن ابتدای خوشه
- به محض پیدا شدن شروع خوشه
- گام‌برداری پیش‌رو تا یافتن مکان اولین باقیمانده برای سطل f_q
- شروع واقعی اجرا آغاز r

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

فیلتر خ

یافتن آغاز خوشه دارای f_q

۱	ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت
۲	$j \leftarrow f_q$
۳	تازمانی که $1 == QF[j].is\ shifted$ انجام بده
۴	$j --$
۵	$r_{\text{آغاز}} \leftarrow j$
۶	تازمانی که $j \neq f_q$
۷	تازمانی که $1 == QF[r_{\text{آغاز}} + 1].runContinued$
۸	$r_{\text{آغاز}} ++$
۹	$r_{\text{آغاز}} ++$
۱۰	تازمانی که $1 \neq QF[j + 1].bucketOccupied$
۱۱	$j ++$
۱۲	$j ++$
۱۳	$r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$
۱۴	تازمانی که $1 == QF[r_{\text{پایان}} + 1].runContinued$ & $1 \leq r_{\text{پایان}} + 1 \leq m$
۱۵	$r_{\text{پایان}} ++$
۱۶	برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

یافتن پایان اجرای دارای f_q

فیلتر خارج قسمت

درج مقداری جدید

- محاسبه خارج قسمت و باقیمانده آن
- در صورت اشغال نبودن سطل متناظر
 - استفاده از روش درج
 - در غیر این صورت، قبل از درج،
 - نیاز به یافتن سطلی مناسب با تابع اسکن از الگوریتم قبلی
 - با یافتن سطل مناسب
 - درج واقعی همچنان نیاز به ادغام مناسب f_r در دنباله عناصر ذخیره شده قبلی
 - امکان نیاز به جابجایی به راست مقادیر بعدی و به روزرسانی بیت‌های متاداده متناظر
- الگوریتم کامل درج با استراتژی انتخاب برای سطل مناسب و تابع جابجایی به راست در الگوریتم ۸

فیلتر خارج قسمت

الگوریتم ۱۰- درج مقدار در فیلتر خارج قسمت

	۱	ورودی: عضو $x \in D$
	۲	ورودی: فیلتر خارج قسمت با تابع درهم‌ساز h
	۳	$f \leftarrow h(x)$
	۴	$f_q, f_r \leftarrow f$
	۵	اگر $QF[f_q].bucketOccupied == 0$ و خالی بودن $QF[f_q]$
خالی و اشغال نبودن سطل	۶	$QF[f_q] \leftarrow f_r$
	۷	$QF[f_q].bucketOccupied \leftarrow 1$
	۸	برگرداندن True
علامت اشغال شدن سطل متناظر	۹	$QF[f_q].bucketOccupied \leftarrow 1$
	۱۰	$r_{\text{پایان}}, r_{\text{تفاز}} \leftarrow Scan(QF, f_q)$
	۱۱	برای i از $r_{\text{تفاز}}$ تا $r_{\text{پایان}}$
به معنای وجود f_r از قبل و لزوم خروج	۱۲	اگر $QF[i] == f_r$
	۱۳	برگرداندن True
درج f_r در سطل i و جابجا کردن بقیه	۱۴	وگر $QF[i] > f_r$ آنگاه
	۱۵	$QF \leftarrow ShiftRight(QF, i)$
	۱۶	$QF[i] \leftarrow f_r$
		اگر $i \neq r_{\text{تفاز}}$
اولین نبودن در اجرا	۱۷	$QF[i].run_continued \leftarrow 1$
	۱۸	اگر $i \neq f_q$
	۱۹	$QF[i].is_shifted \leftarrow 1$
	۲۰	برگرداندن True
باقیمانده جدید بزرگتر از تمامی مقادیر	۲۱	$QF \leftarrow ShiftRight(QF, r_{\text{پایان}} + 1)$
فعلی در اجرا، درج در انتهای اجرا و بسط	۲۲	$QF[r_{\text{پایان}} + 1] \leftarrow f_r$
همان اجرا	۲۳	$QF[i].run_continued \leftarrow 1$
	۲۴	اگر $r_{\text{پایان}} + 1 \neq f_q$
	۲۵	$QF[r_{\text{پایان}} + 1].is_shifted \leftarrow 1$
	۲۶	برگرداندن True

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که $1 == QF[j].isShifted$ انجام یده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که $1 == QF[r_{\text{آغاز}} + 1].runContinued$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که $1 \neq QF[j + 1].bucketOccupied$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که $1 == QF[r_{\text{پایان}} + 1].runContinued$ & $m + 1 \leq r_{\text{پایان}}$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

$$f_q = 12$$

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که ۱ $QF[j].is_shifted ==$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که ۱ $QF[r_{\text{آغاز}} + 1].runContinued ==$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که ۱ $QF[j + 1].bucketOccupied \neq$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که ۱ $r_{\text{پایان}} + 1 \leq m \ \& \ QF[r_{\text{پایان}} + 1].runContinued ==$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

گام ۱

- پیدا کردن شروع خوشه
- گذر به چپ تا زمان وجود انتقال

$$j = f_q = 12$$

اندیس ۱۲: $is_shifted = 1$ پس $j = 11$

اندیس ۱۱: $is_shifted = 1$ پس $j = 10$

اندیس ۱۰: $is_shifted = 0$ پس توقف

$$r_{start} = j = 10$$

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که ۱ $QF[j].is\ shifted ==$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ $QF[r_{\text{آغاز}} + 1].runContinued == 1$ تازمانی که ۱

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که ۱ $QF[j + 1].bucketOccupied \neq$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که ۱ $r_{\text{پایان}} + 1 \leq m \ \& \ QF[r_{\text{پایان}} + 1].runContinued ==$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

گام ۲

- گذر از اجراها تا رسیدن j به $f_q = 12$
- تکرار اول ($j = 10, r_{start} = 10, 12 \neq 10$)
- زیرگام ۲ الف
- یافتن r_{start}
- $r_{start} = 10$
- $runContinued = 1$ بررسی $r_{start} + 1 = 11$
- پس $r_{start} = 11$
- $runContinued = 1$ بررسی $r_{start} + 1 = 12$
- پس $r_{start} = 12$
- $runContinued = 0$ بررسی $r_{start} + 1 = 13$
- پس توقف
- $r_{start} = 13$

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که $1 == QF[j].is\ shifted$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که $1 == QF[r_{\text{آغاز}} + 1].runContinued$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که $1 \neq QF[j + 1].bucketOccupied$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که $1 == QF[r_{\text{پایان}} + 1].runContinued$ & $r_{\text{پایان}} + 1 \leq m$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

زیرگام ۲ ب:

انتقال j به سطل اشغال شده بعدی

▪ $j = 10$

▪ $i+1 = 11$ بررسی $bucket_occupied=1$

▪ پس توقف

$j = 11$

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که ۱ $QF[j].is\ shifted ==$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که ۱ $QF[r_{\text{آغاز}} + 1].runContinued ==$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که ۱ $QF[j + 1].bucketOccupied \neq$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که ۱ $r_{\text{پایان}} + 1 \leq m \ \& \ QF[r_{\text{پایان}} + 1].runContinued ==$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

پایان تکرار اول

$j = 11$

$r_{start} = 13$

$12 \neq 11$ پس ادامه

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که ۱ $QF[j].is\ shifted ==$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ $QF[r_{\text{آغاز}} + 1].runContinued == 1$ تازمانی که ۱

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که ۱ $QF[j + 1].bucketOccupied \neq$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که ۱ $r_{\text{پایان}} + 1 \leq m \ \& \ QF[r_{\text{پایان}} + 1].runContinued ==$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

▪ تکرار دوم ($r_{\text{start}} = 13, j = 11$)

▪ زیرگام ۲ الف

▪ $r_{\text{start}} = 13$

▪ $runContinued = 1$ بررسی $r_{\text{start}} + 1 = 14$

▪ پس $r_{\text{start}} = 14$

▪ $runContinued = 1$ بررسی $r_{\text{start}} + 1 = 15$

▪ پس توقف

▪ $r_{\text{start}} = 15$

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که ۱ $QF[j].is\ shifted ==$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که ۱ $QF[r_{\text{آغاز}} + 1].runContinued ==$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که ۱ $QF[j + 1].bucketOccupied \neq$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که ۱ $r_{\text{پایان}} + 1 \leq m \ \& \ QF[r_{\text{پایان}} + 1].runContinued ==$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

زیرگام ۲ ب:

▪ $j = 11$

▪ $i+1 = 12$ بررسی $bucket_occupied=1$

▪ پس توقف

▪ $j = 12$

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که ۱ $QF[j].is\ shifted ==$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که ۱ $QF[r_{\text{آغاز}} + 1].runContinued ==$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که ۱ $QF[j + 1].bucketOccupied \neq$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که ۱ $r_{\text{پایان}} + 1 \leq m \ \& \ QF[r_{\text{پایان}} + 1].runContinued ==$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

پایان تکرار دوم

$j = 12$

$r_{start} = 15$

$f_q = 12$

▪ خارج شدن از حلقه

الگوریتم ۹- پیمایش فیلتر خارج قسمت برای یافتن اجرا

f_q	اشغال	ادامه	انتقال
10	1	0	0
11	1	1	1
12	1	1	1
13	1	0	1
14	1	1	1
15	0	0	1
16	0	1	1
17	0	0	1
18	0	0	1
19	0	0	0

۱ ورودی: اندیس سطل متعارف f_q ، فیلتر خارج قسمت

۲ $j \leftarrow f_q$

۳ تازمانی که $1 == QF[j].is\ shifted$ انجام بده

۴ $j --$

۵ $r_{\text{آغاز}} \leftarrow j$

۶ تازمانی که $j \neq f_q$

۷ تازمانی که $1 == QF[r_{\text{آغاز}} + 1].runContinued$

۸ $r_{\text{آغاز}} ++$

۹ $r_{\text{آغاز}} ++$

۱۰ تازمانی که $1 \neq QF[j + 1].bucketOccupied$

۱۱ $j ++$

۱۲ $j ++$

۱۳ $r_{\text{آغاز}} \leftarrow r_{\text{پایان}}$

۱۴ تازمانی که $1 == QF[r_{\text{پایان}} + 1].runContinued$ و $1 \leq r_{\text{پایان}} + 1 \leq m$

۱۵ $r_{\text{پایان}} ++$

۱۶ برگرداندن $r_{\text{آغاز}}$ و $r_{\text{پایان}}$

فیلتر خارج قسمت

گام ۳

- یافتن پایان اجرای شروع شده از ۱۵
- $r_{\text{end}} = r_{\text{start}} = 15$
- $r_{\text{end}} + 1 = 16$ بررسی $runContinued = 1$
- پس $r_{\text{end}} = 16$
- $r_{\text{end}} + 1 = 17$ بررسی $runContinued = 0$
- پس توقف
- خروجی: $(r_{\text{end}}, r_{\text{start}}) = (16, 15)$

فیلتر خارج قسمت

با توجه به طرح کاوش خطی

▪ زمان اکثر اجراها از مرتبه $O(1)$

▪ اشاره نویسندگان فیلتر به احتمال زیاد همه اجراها دارای طول $O(\log m)$

فیلتر خارج قسمت

مثال - افزودن مقادیر به فیلتر

فیلتر خارج قسمت با اثر انگشت‌های ۱۶ بیتی حاصل از نسخه ۳۲ بیتی تابع درهم‌ساز مورمور^۳

$$h(x) = \text{MurmurHash3}(x) \% 16$$

سطل‌بندی، رزرو $q = 3$ بیت با بیشترین اهمیت

▪ تبعاً اندازه فیلتر باقیمانده برابر $m = 2^3 = 8$

▪ ذخیره $p = 13$ بیت باقیمانده در سطل‌های انتخاب شده

فیلتر خارج قسمت

همانند مثال قبلی نمایه‌سازی نام پایتخت‌ها

کپنهاگ اولین مقدار افزودنی به فیلتر

محاسبه اثر انگشت آن با استفاده از تابع درهم‌ساز h

$$f = h(\text{Copenhagen}) = 4248224207$$

▪ خارج قسمت و باقیمانده

$$f_q = \left\lfloor \frac{f}{2^{13}} \right\rfloor = 7, f_r = f \% 2^{13} = 490127823$$

فیلتر خارج قسمت

همانند مثال قبلی نمایه‌سازی نام پایتخت‌ها

کپنهاگ اولین مقدار افزودنی به فیلتر

محاسبه اثر انگشت آن با استفاده از تابع درهم‌ساز h

$$f = h(\text{Copenhagen}) = 4248224207$$

خارج‌قسمت و باقیمانده

$$f_q = \left\lfloor \frac{f}{2^{13}} \right\rfloor = 7, f_r = f \% 2^{13} = 490127823$$

0	1	2	3	4	5	6	7														
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
																					490127823

سطل متناظر کپنهاگ $j = f_q = 7$

پر کردن باقیمانده f_r در آن

درج در این نقطه ساده به دلیل پر نبودن سطل‌ها

درج $f_r = 490127823$ در سطل با نمایه ۷ و قرار دادن مقدار یک در بیت `bucket_occupied`

فیلتر خارج قسمت

مقدار لیسبون

▪ اثر انگشت $f = 629555247$ و سطل متناظر ۱،

مقدار پاریس

▪ اثر انگشت $f = 2673248856$ و سطل متناظر ۴

به دلیل آزاد بودن سطل‌های متناظر، درج به ترتیب باقیمانده‌ها و تنظیم بیت‌های `bucket_occupied`

0	1	2	3	4	5	6	7													
0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
			92684335						525765208									490127823		

فیلتر خارج قسمت

درج استکھلم

- اثر انگشت $f = 775943400$ و سطل متناظر $j = f_q = 1$ و باقیمانده $f_r = 239072488$
- مقداردهی بیت `bucket_occupied` سطل متناظر ۱
- به معنای اشغال با باقیمانده مقداری دیگری (در این مورد لیسبون)

فیلتر خارج قسمت

درج زاگرب

اثر انگشت $f = 1474643542$ ، سطل متناظر $j = 2$ و باقیمانده $f_r = 400901718$

بیت مقداره‌ی شده `is_shifted` نمایشگر اشغال سطل ۲ با مقداری انتقالی

فیلتر خارج قسمت

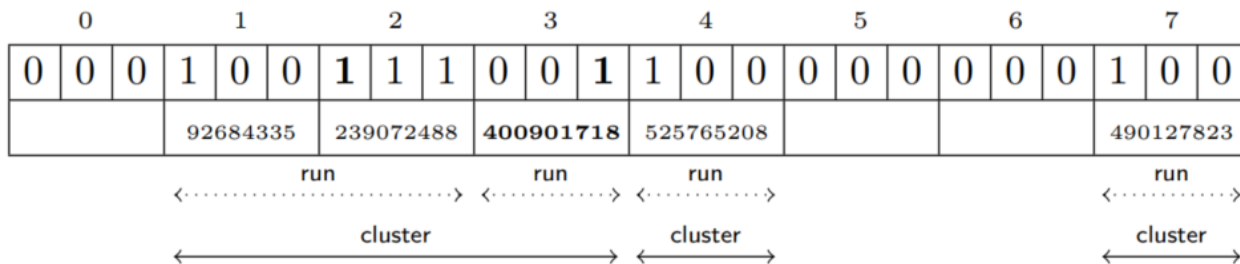
درج زاگرب

اثر انگشت $f = 1474643542$ ، سطل متناظر $j = 2$ و باقیمانده $f_r = 400901718$

بیت مقداردهی شده `is_shifted` نمایشگر اشغال سطل ۲ با مقداری انتقالی

عدم تنظیم بیت `bucket_occupied`

▪ انتقال مقدار f_r به سمت راست، به سطل بعدی موجود که در این مورد سطل ۳



فیلتر خارج قسمت

مقداردهی بیت `run_shifted`

- نمایش سطل حاوی مقداری منتقل شده از موقعیت متناظر خود
- عدم تغییر بیت `run_shifted` به دلیل متناظر بودن اولین عضو آن با سطل مرتبط
- مضافاً، مقداردهی بیت `bucket_occupied` سطل ۲ جهت به یاد داشتن ذخیره شدن حداقل یک باقیمانده مربوط سطل متناظر

درج ورشو

▪ با اثر انگشت $f = 567538184$ با خارج قسمت و باقیمانده

$$f_q = \left\lfloor \frac{f}{2^{13}} \right\rfloor = 1$$

$$f_r = f \% 2^{13} = 30667272$$

فیلتر خارج قسمت

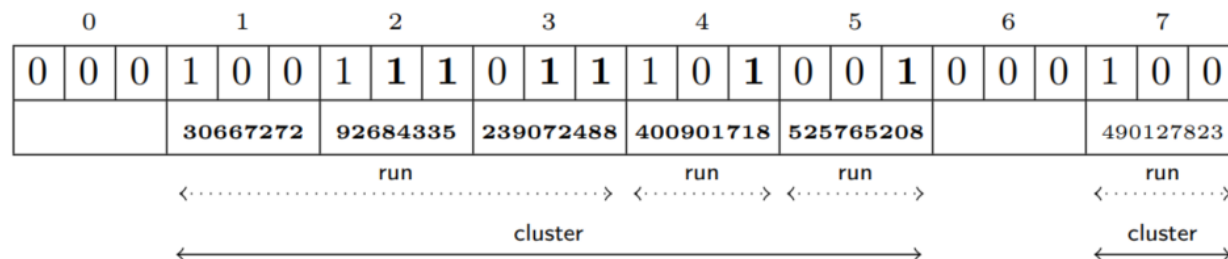
اشغال از قبل سطل متناظر $j = f_q = 1$ بر اساس بیت `bucket_occupied`

- عدم مقداردهی بیت‌های دیگر
- به معنای در ابتدای خوشه بودن
- کوچکتر بودن مقدار باقیمانده f_q از مقدار نمایه شده ۹۲۶۸۴۳۳۵
- در نتیجه نیاز به نمایه شدن در سطل متناظر و انتقال بقیه باقیمانده‌ها و مقداردهی به عنوان ادامه

فیلتر خارج قسمت

اشغال از قبل سطل متناظر $j = f_q = 1$ بر اساس بیت bucket_occupied

- عدم مقداردهی بیت‌های دیگر
- به معنای در ابتدای خوشه بودن
- کوچکتر بودن مقدار باقیمانده f_q از مقدار نمایه شده ۹۲۶۸۴۳۳۵
- در نتیجه نیاز به نمایه شدن در سطل متناظر و انتقال بقیه باقیمانده‌ها و مقداردهی به عنوان ادامه



فیلتر خارج قسمت

آزمون عضویت

- همانند درج انجام
- بررسی سطل متناظر
- مشاهده وجود باقیمانده متناظر در فیلتر با مشاهده بیت `bucket_occupied`
- در صورت یک نبودن مقدار بیت مذکور،
- عدم وجود مقدار در فیلتر
- وگرنه، بررسی فیلتر با رویه اسکن جهت یافتن اجرای مناسب سطل
- مقایسه باقیمانده‌های حاضر در اجرا با باقیمانده مقدار مورد پرسش
- در صورت یافت شدن برگرداندن احتمال وجود آن در فیلتر

فیلتر

الگوریتم ۱۱- آزمون عنصر در فیلتر خارج قسمت

۱ ورودی: عنصر $x \in D$

۲ ورودی: فیلتر خارج قسمت با تابع درهم ساز h

۳ خروجی: False اگر عنصر یافت نشد و True اگر احتمال حضور وجود داشته باشد

۴ $f \leftarrow h(x)$

۵ $f_q, f_r \leftarrow f$

۶ اگر $1 \neq QF[f_q].bucket_occupied$ آنگاه

۷ برگرداندن False

۸ در غیر این صورت

۹ $r_{\text{پایان}}, r_{\text{آغاز}} \leftarrow \text{Scan}(QF, f_q)$

۱۰ برای i از $r_{\text{آغاز}}$ تا $r_{\text{پایان}}$ انجام بده

۱۱ اگر $QF[i] == f_r$ آنگاه

۱۲ برگرداندن True

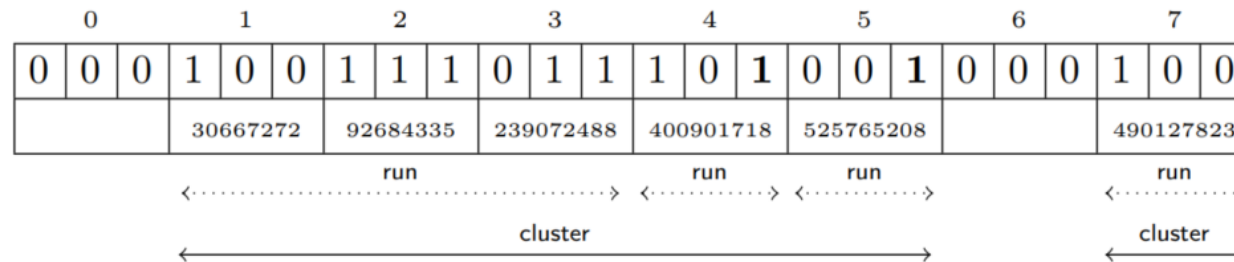
۱۳ برگرداندن False

جستجوی f_r در داخل اجرا

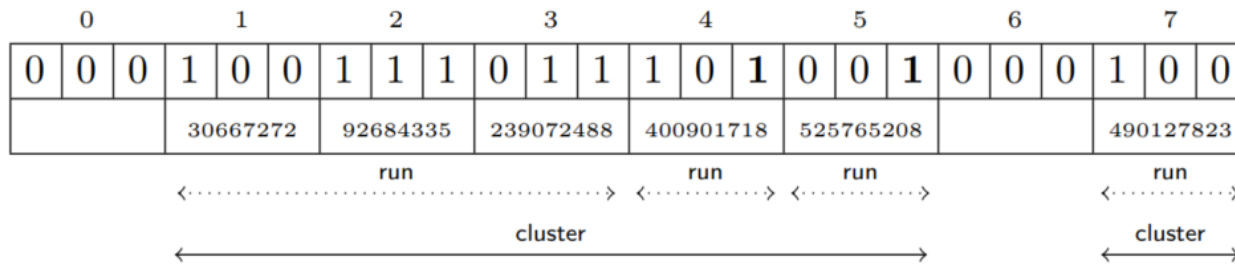
فیلتر خارج قسمت

مثال - آزمون عناصر در فیلتر

داده ساختار حاصل مثال قبل



فیلتر خارج قسمت



آزمون عضویت پاریس با خارج قسمت $f_q = 4$ و باقیمانده $f_r = 525765208$

- اشغال سطل ۴ از قبل

- به معنای وجود حداقل یک باقیمانده در جایی از فیلتر با سطل متناظر ۴
- عدم امکان مقایسه مقدار سطل را با f_r به دلیل تنظیم بیت `bucket_shifted`
- در نتیجه در پی اجرایی دارای تطابق با سطل متناظر ۴ در خوشه فعلی
- اسکن از سطل ۴ به سمت چپ
- شمارش سطل‌ها را با مقدار بیت‌های `bucket_occupied` جهت رسیدن به آغاز خوشه
- آغاز خوشه از سطل ۱ و دو سطل اشغال شده (سطل‌های ۱ و ۲) در سمت چپ سطل ۴ قرار دارند. بنابراین، اجرای مدنظر سومین در خوشه
- در نتیجه نیاز به اسکن از ابتدای خوشه (سطل ۱) به سمت راست و شمارش سطل‌هایی که بیت‌های `run_shifted` برابر صفر تا رسیدن به اجرای مدنظر
- آغاز اجرا از سطل ۵
- مقایسه مقدار با باقیمانده‌های ذخیره
- با در نظر گرفتن ترتیب صعودی آنها

مطابقت مقدار موجود در سطل ۵ با باقیمانده $f_r = 525765208$

- امکان وجود پاریس در فیلتر

فیلتر خارج قسمت

حذف از فیلتر خارج قسمت

به روشی سخت شبیه به درج مقدار جدید

ذخیره با ترتیب افزایشی باقیمانده‌های اثر انگشت‌ها با خارج قسمت یکسان

حذف باقیمانده‌ای از خوشه نیاز به جابجایی تمام اثر انگشت‌ها و اصلاح بیت‌های به ترتیب متاداده

فیلتر خارج قسمت

الگوریتم ۱۲- استفاده از انتقال به چپ برای پر کردن سطل‌های خالی

۱	ورودی: اندیس سطل i
۲	ورودی: فیلتر خارج قسمت به طول m
۳	$k \leftarrow i + 1$
۴	تازمانی که $QF[k] \neq \emptyset$ انجام بده
۵	$QF[k - 1] \leftarrow QF[k]$
۶	$QF[k - 1].is_shifted \leftarrow QF[k].is_shifted$
۷	$QF[k] \leftarrow \phi$
۸	$QF[k].run_continued \leftarrow \cdot$
۹	$QF[k].is_shifted \leftarrow \cdot$
۱۰	$k \leftarrow k + 1$
۱۱	اگر $k > m$ آنگاه
۱۲	$k \leftarrow \cdot$

فیلتر خارج قسمت

نیاز به بررسی اشغال سطل متناظر

▪ در صورت عدم مقدارهی عدم وجود قطعی در فیلتر و توقف

در صورت مقدارهی

▪ اسکن برای یافتن سطل مناسب و حذف عنصر درخواستی (در صورت وجود) و جابجائی مقادیر عناصر بعدی و به روزرسانی بیت‌های متاداده متناظر

در صورت حذف آخرین باقیمانده سطل متناظر پاکسازی بیت `bucket_occupied`

الگوریتم ۱۳ - حذف مقدار از فیلتر خارج قسمت

	۱	ورودی: عنصر $x \in D$
	۲	ورودی: فیلتر خارج قسمت با تابع درهم ساز h
	۳	خروجی: False اگر عنصر یافت نشد و True در غیر این صورت
	۴	$f \leftarrow h(x)$
	۵	$f_q, f_r \leftarrow f$
	۶	اگر $1 \neq QF[f_q].bucket_occupied$ آنگاه برگرداندن True
	۷	$r_{\text{پایان}}, r_{\text{آغاز}} \leftarrow Scan(QF, f_q)$
	۸	برای i از آغاز r تا پایان r انجام بده
	۹	اگر $QF[i] == f_r$ آنگاه
یافتن مقدار و حذف آن	۱۰	$QF[i] \leftarrow \phi$
	۱۱	$QF[i].run_continued \leftarrow \cdot$
	۱۲	$QF[i].isShifted \leftarrow \cdot$
انتقال به چپ	۱۳	اگر پایان $r < i$ آنگاه
	۱۴	$QF \leftarrow ShiftLeft(QF, i)$
تهی شدن اجرا	۱۵	اگر پایان $r == r_{\text{آغاز}}$ آنگاه
صفر شدن مقدار اشغال	۱۶	$QF[i].bucket_occupied \leftarrow \cdot$
	۱۷	برگرداندن True
نبودن عضو در اجرا	۱۸	برگرداندن False

فیلتر خارج

فیلتر خارج قسمت-ویژگی ها

مثبت کاذب امکان پذیر

- داده ساختار فیلتر خارج قسمت نمایش فشرده ای از مجموعه **چندگانه** از اثر انگشت ها
- نسبت مثبت کاذب آن تابعی از تابع درهم ساز h و تعداد مقادیر عناصر n اضافه شده به فیلتر
- مضافاً، تصادم سخت
- دو مقدار متفاوت دارای مقادیر یکسان باقیمانده و خارج قسمت
- به دلیل چنین رویدادهای بسیار نادری، امکان وقوع پاسخهای مثبت کاذب و محدود بودن احتمال آنها pr_{fp} با کران بالای زیر

$$pr_{fp} \approx 1 - e^{-\frac{n}{2^p}} \leq \frac{n}{2^p}$$

پس، با تعداد ثابت عناصر مورد انتظار n ، امکان سبک-سنگینی بین احتمال مثبت کاذب pr_{fp} و طول اثر انگشت p

استفاده پیاده سازی های عملی فیلترهای خارج قسمت از اثر انگشت های ۳۲ و ۶۴ بیتی

فیلتر خارج قسمت

اهمیت ضریب بار در فیلتر خارج قسمت همانند سایر جدول های درهم

در پی تخصیص حداقل اندازه به تعداد مقادیر مورد انتظار سطل

▪ به بیان دیگر، انتخاب تعداد سطل ها m به صورت زیر

$$m = 2^q > n$$

محاسبه طول باقیمانده r

$$r = \left\lceil \log \left(-\frac{n}{2^q} \times \frac{1}{\ln(1-prfp)} \right) \right\rceil$$

فیلتر خارج قسمت

منفی کاذب امکان ناپذیر

▪ در صورت پاسخ منفی الگوریتم عضویت فیلتر خارج قسمت، قطعاً عدم عضویت در مجموعه

$$pr_{fn} = 0$$

سرعت

▪ فیلتر خارج قسمت حدود بیست درصد بزرگتر از فیلتر بلوم

▪ در عین حال سریع تر

▪ چرا؟

▪ سرعت بالاتر به دلیل نیاز به ارزیابی یک تابع درهم ساز در هر دسترسی

▪ ذخیره داده ها در **پیمانه های پیوسته**

▪ آزمون ها در فیلتر خارج قسمت منجر به یک **missed cache**

▪ در مقابل حداقل دو **missed cache** در انتظار برای الگوریتم فیلتر بلوم

فیلتر خارج قسمت

میزان مثبت کاذب و ملاحظات فضایی

- وقوع مثبت کاذب در فیلتر خارج قسمت با تولید اثرانگشت یکسان برای دو کلید متفاوت
- جدول شامل 2^q خانه و طول اثرانگشت برابر $p = q + r$
- آنگاه احتمال وقوع مثبت کاذب تقریباً برابر با 2^{-r}
- وابستگی این احتمال به تعداد بیت‌های باقی‌مانده
- فضای موردنیاز برای جدول فیلتر خارج قسمت برابر با $2^q(r + 3)$ بیت
- تعداد عناصری که در فیلتر درج می‌شوند برابر با $n = \alpha 2^q$
- α ضریب بارگذاری
- تاثیر چشمگیر ضریب تأثیر بر کارایی عملیات درج، جست‌وجو، و حذف

معرفی گونه کم‌حجم‌تر از فیلتر خارج قسمت

- استفاده از تنها از دو بیت متاداده
- نیاز به فضای $2^q(r + 2)$ بیت
- با وجود کاهش یک بیت، عدم افزایش میزان مثبت کاذب
- پیچیدگی سخت بیشتر مرحله رمزگشایی
- به‌ویژه در مواقع وجود خوشه‌های بلند،
- عملیات متداول (درج، جست‌وجو و حذف) وابستگی بیشتر به پردازنده و پرهزینه‌تر

فیلتر خارج قسمت

مقایسه با فیلترهای بلوم

- به سبب فضای اضافی لازم برای جدول جست و جوی خطی
- برای میزان رایج مثبت کاذب فیلترهای خارج قسمت معمولاً مصرف فضا اندکی بیشتر از فیلترهای بلوم فضا
- کارایی فضایی بیشتر فیلترهای خارج قسمت نسبت به فیلترهای بلوم برای میزان بسیار پایین مثبت کاذب
- درج
- فیلترهای خارج قسمت می توانند ۲,۴ میلیون درج در ثانیه در مقابل فیلترهای بلوم محدود به حدود ۰,۶۹ میلیون
- آزمون عضویت
- هر دو در همان سطح حدود ۲ میلیون در ثانیه

فیلتر خارج قسمت

مثال - حافظه مورد نیاز -

- برای مدیریت یک میلیارد مقدار، فیلتر خارج قسمت باید حداقل 2^{30} سطل
- عدم امکان استفاده از اثر انگشت‌های کوتاه‌تر از ۳۱ بیت
- در پی احتمال رویدادهای مثبت کاذب حدود دو درصد
- تعداد بیت‌ها برای باقیمانده

$$r = \left\lceil \log \left(-\frac{10^9}{2^{30}} \times \frac{1}{\ln(1 - 0.02)} \right) \right\rceil = 6$$

- پس، طول مورد نیاز اثر انگشت $p = q + r = 36$ بیت
- استفاده از ۳۰ بیت اول برای سطل‌بندی و ذخیره ۶ بیت باقی‌مانده در سطل مناسب
- هر سطل علاوه بر این شامل سه بیت متاداده
- اندازه کل فیلتر خارج قسمت 9×2^{30} بیت معادل تقریباً ۱,۲ گیگابایت حافظه

فیلتر خارج قسمت

فیلتر خارج قسمت

امکان بازیابی اثر انگشتها از دادههای ذخیره شده

- در نتیجه پشتیبانی از حذف، ادغام، و تغییر اندازه

عدم تاثیر ادغام بر نسبت مثبت کاذب فیلترها

همچنین،

- پشتیبانی فیلتر بلوم شمارنده از حذفهای احتمالا درست اما درستی همیشگی حذف در فیلتر خارج قسمت

فیلتر خارج قسمت

تغییر اندازه فیلتر خارج قسمت

▪ هم کوچک سازی و هم بزرگ سازی

▪ با حرکت روی فیلتر

▪ کپی کردن هر اثر انگشت در داده ساختاری تازه تخصیص یافته

▪ بدون نیاز به درهم سازی مجدد

خق	باقیمانده	متاداده		
۰۰۰۰				
۰۰۰۱				
۰۰۱۰				
۰۰۱۱				
۰۱۰۰	۱۰۱۰	۱	۰	۰
۰۱۰۱				
۰۱۱۰				
۰۱۱۱				
۱۰۰۰				
۱۰۰۱				
۱۰۱۰				
۱۰۱۱				
۱۱۰۰				
۱۱۰۱				
۱۱۱۰	۰۱۱۱	۱	۱	۰
۱۱۱۱				

خق	باقیمانده	متاداده		
۰۰۰				
۰۰۱				
۰۱۰	۰۱۰۱۰	۱	۰	۰
۰۱۱				
۱۰۰				
۱۰۱				
۱۱۰				
۱۱۱	۰۰۱۱۱	۱	۰	۰

خق	باقیمانده	متاداده		
۰۰				
۰۱	۰۰۱۰۱۰	۱	۰	۰
۱۰				
۱۱	۱۰۰۱۱۱	۱	۰	۰

فیلتر خارج قسمت

ادغام دو یا چند فیلتر خارج قسمت

▪ با استفاده از الگوریتمی شبیه به مرتب‌سازی ادغامی، الگوریتم مرتب‌سازی تقسیم و غلبه ابداعی جان فون نویمان

امکان اسکن موازی همه فیلترهای ورودی و نوشتن نتیجه ادغام شده در فیلتر خروجی

زمان مورد نیاز برای انجام آزمون، افزودن یا حذف در فیلتر خارج قسمت

▪ وابسته به زمان اسکن به عقب و جلو

طراحی فیلتر خارج قسمت با تمرکز بر کلان‌داده (به عنوان مثال، یک میلیارد مقدار برای تابع درهم‌ساز ۶۴ بیتی)

کاهش مزایا در مجموعه‌های داده کوچک یا متوسط، به دلیل پیچیدگی آن

مرتب‌سازی ادغامی

اگر $n = 1$ ، آن‌گاه آرایه مرتب

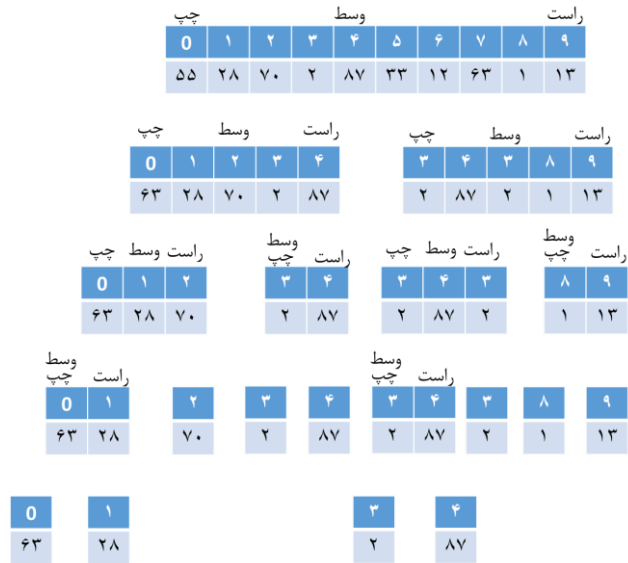
در صورت مرتب بودن جداگانه دو زیرفهرست با اندیس‌های $[0:n/2]$ و $[n/2 + 1:n - 1]$

- ادغام روشمند موجب مرتب‌سازی ترکیب دو زیرفهرست
- دو عمل تقسیم و ترکیب
- عدم جابجایی در هنگام تقسیم

مرتب‌سازی ادغامی

چپ	وسط							راست	
0	1	2	3	4	5	6	7	8	9
55	28	70	2	87	33	12	63	1	13

مرتب‌سازی ادغامی



مرتب‌سازی ادغامی

چپ	وسط					راست			
0	1	2	3	4	5	6	7	8	9
55	28	70	2	87	33	12	63	1	13

چپ	وسط		راست		چپ	وسط		راست	
0	1	2	3	4	3	4	3	8	9
63	28	70	2	87	2	87	2	1	13

چپ	وسط		راست		وسط	راست		چپ	وسط		راست	
0	1	2	3	4	3	4	3	8	9	1	13	
63	28	70	2	87	2	87	2	1	13			

وسط	راست		وسط		چپ		راست		وسط		راست	
0	1	2	3	4	3	4	3	8	9	1	13	
63	28	70	2	87	2	87	2	1	13			

0	1	3	4
63	28	2	87

0	1	5	6
63	28	33	12

0	1	2	3	4	5	6	7	8	9
28	63	70	2	87	12	33	55	1	13

0	1	2	3	4	5	6	7	8	9
28	63	70	2	87	12	33	55	1	13

0	1	2	3	4	5	6	7	8	9
2	28	63	70	87	1	12	13	33	55

0	1	2	3	4	5	6	7	8	9
1	2	12	13	28	33	55	63	70	80

مرتب‌سازی ادغام

شبه‌کد بخش ادغام آن به صورت زیر است.

```
Alg. Edqam(arr, c, v, r)
  arr_c = arr[c:v]; arr_r[v+1:r]
  i = j = c
  k = c
  while i <= v & j <= r - v
    if arr_c[i] <= arr_r[j]
      arr[k] = arr_c[i++]
    else
      arr[k] = arr_r[j++]
    k++
  while i <= v
    arr[k] = arr_c[i++]
    i++; k++
  while j <= r - v
    arr[k] = arr_r[j++]
    j++; k++
```

الگوریتم تقسیم و غلبه بر اساس تابع ادغام به صورت زیر تعریف می‌شود.

```
Alg. Morattabsazi_edqami(arr, c, r)
  if c < r
    v = c + (r - c) / 2
    Morattabsazi_edqami(arr, c, v)
    Morattabsazi_edqami(arr, v + 1, r)
  Edqam(arr, c, v, r)
```

فیلتر فاخته

اکثر بهبودهای فیلتر بلوم کلاسیک با پشتیبانی از حذف

دچار افت از نظر فضا یا عملکرد

بین فن، دیوید اندرسون، مایکل کامینسکی و مایکل میتزن ماخر در سال ۱۳۹۳ پیشنهاد فیلتر فاخته

نوعی فشرده از جدول درهم فاخته

به جای ذخیره جفت‌های کلید-مقدار، تطبیق اثر انگشت‌هایی با طول p برای هر مقدار درج شده

پیاده‌سازی آسان‌تر، پشتیبان درج و حذف پویا

استفاده از فضای کمتری نسبت به سایر اصلاحات فیلتر بلوم در بسیاری از کاربردهای عملی

نمایش عملکردی بالاتر

فیلتر فاخته

داده‌ساختار با جدول درهم فاخته چندگانه با m سطل

ذخیره b مقدار در هر یک

درج مقدار جدید در درهم فاخته استاندارد

- نیاز به دسترسی به مقادیر موجود اصلی
- جهت تعیین محل جدید برای ذخیره مقادیر جدید در صورت نیاز به فضا
- ذخیره اثر انگشت‌ها در فیلتر فاخته
- عدم وجود راهی برای بازیابی عناصر اصلی و درهم‌سازی مجدد آنها برای یافتن سطل جدیدشان در جدول درهم

فیلتر فاخته

هدف غلبه بر این محدودیت و همچنان استفاده از درهم‌ساز فاخته

- استفاده الگوریتم فیلتر فاخته از «درهم‌ساز فاخته با کلید جزئی»
- استخراج مقدار موجود از سطل بدون دانستن خود مقدار اصلی و صرفاً از اثر انگشت آن
- درج هر مقدار x ، محاسبه اثر انگشت p -بیتی f و نمایه‌های دو سطل نامزد:

$$i = h(x)\%m$$

$$j = (i \oplus (h(f)\%m))\%m$$

جهت داشتن توزیعی یکنواخت از مقادیر

فیلتر فاخته

کوچکتر بودن طول اثر انگشت p در مقایسه با طول فیلتر m

- تغییر تعداد کمی از بیت‌های پایین‌تر با عملیات یای انحصاری
- ثابت ماندن بیشتر بیت‌های مرتبه بالاتر
- نمایشگر تمایل مقادیر انتقال یافته از سطل‌های اصلی خود به قرارگیری در نزدیکی یکدیگر در سطل‌های جانشین
- در نتیجه توزیع در جدول درهم دچار انحراف
- تاثیرگذار بر کارایی فیلتر
- درهم‌ساز اثر انگشت‌ها ضامن انتقال این مقادیر به سطل‌هایی در بخش‌های کاملاً متفاوت جدول درهم
- بنابراین کاهش تصادم‌های درهم و بهبود استفاده از جدول

فیلتر فاخته

عملیات یا انحصاری \oplus

▪ با دانستن سطل‌های فعلی مقدار k امکان محاسبه سطل جانشین آن k^* بدون بازیابی مقدار اصلی

$$k^* = (k \oplus h(f)) \% m$$

فیلتر فاخته

درج مقدار جدید x

- محاسبه نمایه‌های دو سطل نامزد با معادلات
- در صورت خالی بودن یکی از دو سطل حاصل، درج مقدار در آن
- در صورت پر بودن هر دو
- انتخاب تصادفی یکی از آن سطل‌ها و ذخیره مقدار x در آن جا
- انتقال مقدار از آن سطل به سطل نامزد جانشین
- تکرار روش تا زمان یافتن سطل خالی یا تا زمان رسیدن به آستانه به حداکثر تعداد جابجایی‌ها
- در صورت خالی نبودن سطل خالی، پر در نظر گرفتن فیلتر

فیلتر فاخته

الگوریتم ۱۴- درج مقدار در فیلتر فاخته

ورودی: مقدار $x \in D$

ورودی: فیلتر فاخته با اثر انگشت گذاری و تابع درهم ساز h

خروجی: True اگر مقدار اضافه شده باشد و False در غیر این صورت

$f \leftarrow$ انگشت اثر (x)

$i \leftarrow h(x)$

$j \leftarrow i \oplus h(f)$

اگر $CUCKOOFILTER [i]$ فضای خالی دارد آنگاه

$CUCKOOFILTER [i].add (f)$

برگرداندن True

در غیر این صورت اگر $CUCKOOFILTER [j]$ فضای خالی دارد آنگاه

$CUCKOOFILTER [j].add (f)$

برگرداندن True

$k \leftarrow$ نمونه $(\{i, j\})$

برای n از صفر تا بیش-تکرار انجام بده

$x \leftarrow$ نمونه $(CUCKOOFILTER [k])$

f و اثر انگشت ذخیره شده در ورودی x را با هم عوض کن

$k = k \oplus h(f)$

اگر $CUCKOOFILTER [k]$ فضای خالی دارد آنگاه

$CUCKOOFILTER [k].add (f)$

برگرداندن True

برگرداندن False

فیلتر فاخته

مثال - افزودن مقدار به فیلتر

داده ساختار به طول $m = 10$

- جهت سادگی، ذخیره صرفاً یک اثر انگشت $p = 16$ بیتی در هر سطل
- استفاده از تابع درهم‌ساز مورمور ۳ با اندازه ۳۲ بیتی برای محاسبه اثر انگشت‌ها و نمایه‌های سطل
 - درج نام پایتخت‌ها درج

فیلتر فاخته

آتن

▪ اثر انگشت $f = 27356$ و سطل‌های نامزد \cdot و γ

▪ سطل اصلی \cdot آزاد

0	1	2	3	4	5	6	7	8	9
27356							49615		

فیلتر فاخته

لیسبون

- اثر انگشت آن $f = 16431$ و سطل‌های نامزد ۷ و ۹
- شروع با سطل اصلی ۷
- اما اشغال شده
- در حداکثر ظرفیت یک
- بررسی سطل جانشین ۹
- خالی
- ذخیره اثر انگشت در آنجا

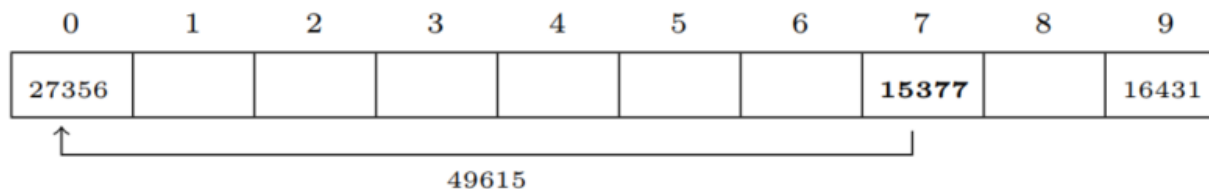
0	1	2	3	4	5	6	7	8	9
27356							49615		16431

فیلتر فاخته

هلسینکی

- اثر انگشت $f = 15377$ و هر دو نمایه سطل برابر با ۷
- چنین برخورد نمایه‌ای برای فیلترهای کوچک، محتمل‌تر از فیلترهای واقعی
- اشغال شدگی سطل ۷ اشغال و عدم امکان پذیرش بیش از یک مقدار
- نیاز به آغاز روش جابجایی در فیلتر
- سطل k
- تبادل مقدار ۴۹۶۱۵ از سطل ۷ با مقدار f
- انتقال آن مقدار به سطل جدید k

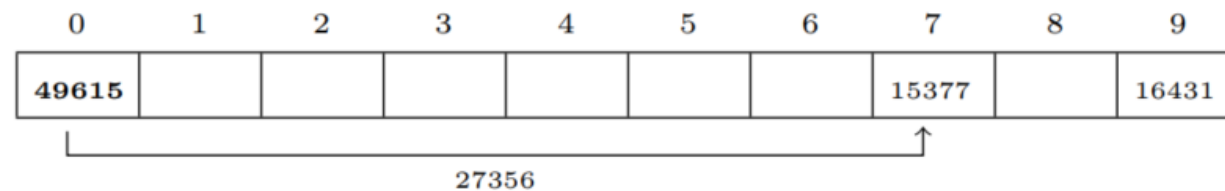
$$k = (7 \oplus \text{MurmurHash3}(49615)) \% 10 = 0$$



فیلتر فاخته

- سطل ۰ قبلاً حاوی مقدار ۲۷۳۵۶
- تبادل آن با ۴۹۶۱۵
- محاسبه نمایه سطل جدید

$$k = (0 \oplus \text{MurmurHash3}(27356)) \% 10 = 7$$



فیلتر فاخته

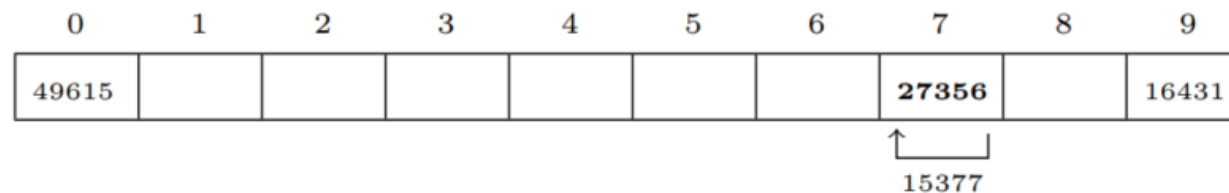
برگشت به سطل ۷
▪ پر بودن آن

نیاز به تکرار روش جابجایی

تبادل مقدار ۲۷۳۵۶ و ذخیره آن در سطل فعلی

محاسبه سطل جدید برای مقدار ۱۵۳۷۷

$$k = (7 \oplus \text{MurmurHash3}(15377)) \% 10 = 7$$



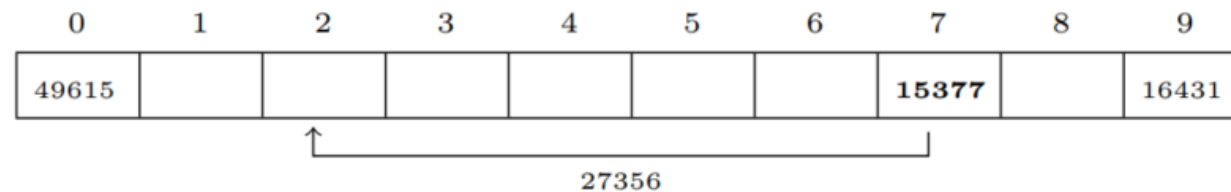
فیلتر فاخته

برگشت به نمایه به دلیل تصادم

ذخیره مقدار ۱۵۳۷۷ در آن

انتقال مقدار ۲۷۳۵۶ به یک سطل جدید k

$$k = (7 \oplus \text{MurmurHash3}(27356)) \% 10 = 2$$



فیلتر فاخته

سطل ۲ خالی

▪ ذخیره مقدار ۲۷۳۵۶ در آن

▪ ختم روش درج

0	1	2	3	4	5	6	7	8	9
49615		27356					15377		16431

فیلتر فاخته

آزمون وجود مقدار در فیلتر

- ساده
- ابتدا، برای عنصر مورد آزمون، محاسبهٔ اثر انگشت و سطل‌های نامزد آن
- اگر اثر انگشت در هر یک از سطل‌ها وجود داشته باشد، نتیجه می‌گیریم که عنصر ممکن است وجود داشته باشد.
- در غیر این صورت، قطعاً در فیلتر نیست.

فیلتر فاخته

الگوریتم ۱۵- آزمون مقدار در فیلتر فاخته

ورودی: مقدار $x \in D$

ورودی: فیلتر فاخته با اثر انگشت‌گذاری و تابع درهم‌ساز h

خروجی: False اگر مقدار یافت نشد و True اگر ممکن است وجود داشته باشد

$f \leftarrow$ انگشت اثر (x)

$i \leftarrow h(x)$

$j \leftarrow i \oplus h(f)$

اگر $f \in \text{CuckooFilter}[i]$ یا $f \in \text{CuckooFilter}[j]$ آنگاه

برگرداندن True

برگرداندن False

فیلتر فاخته

مثال - آزمون عناصر در فیلتر

0	1	2	3	4	5	6	7	8	9
49615		27356					15377		16431

فیلتر فاخته

لیسبون

با سطل‌های نامزد آن ۷ و ۹ و اثر انگشت آن $f = 16431$

همانطور که در بالا محاسبه کردیم، آزمایش کنیم.

می‌توانیم مقدار ۱۶۴۳۱ را در سطل ۹ پیدا کنیم و نتیجه بگیریم که عنصر لیسبون ممکن است در فیلتر وجود داشته باشد.

در مقابل، عنصر اسلو را در نظر می‌گیریم که دارای اثر انگشت $f = 53104$ و سطل‌های نامزد ۰ و ۶ است.

همانطور که می‌بینیم، چنین مقداری در این سطل‌ها وجود ندارد، بنابراین عنصر اسلو قطعاً در فیلتر نیست.

فیلتر فاخته

حذف عنصر،

ایجاد اثر انگشت

محاسبه نمایه‌های سطل‌های نامزد

بررسی اثر انگشت در نمایه‌ها

در صورت مطابقت با هر یک از مقادیر موجود در هر یک از سطل‌ها

▪ حذف یک نسخه از اثر انگشت از آن سطل

فیلتر فاخته

الگوریتم ۱۶- حذف مقدار از فیلتر فاخته

ورودی: مقدار $x \in D$

ورودی: فیلتر فاخته با اثر انگشت گذاری و تابع درهم ساز h

خروجی: True اگر مقدار حذف شده باشد و False در غیر این صورت

$f \leftarrow$ انگشت اثر (x)

$i \leftarrow h(x)$

$j \leftarrow i \oplus h(f)$

اگر $f \in \text{CUCKOOFILTER}[i]$ آنگاه

`CUCKOOFILTER[i].drop(f)`

برگرداندن True

در غیر این صورت اگر $f \in \text{CUCKOOFILTER}[j]$ آنگاه

`CUCKOOFILTER[j].drop(f)`

برگرداندن True

برگرداندن False

فیلتر فاخته-ویژگی‌ها

دو راه حل پشتیبانی از حذف فیلتر فاخته نیاز

- ذخیره چندین نسخه از مقدار

- ایجاد شمارنده برای هر مقدار ذخیره شده

هر دو رویکرد پشتیبانی «حذف‌های احتمالا درست»

- دلایل:

- به دلیل ظرفیت محدود سطل

- عدم امکان ذخیره بیش از $2b$ مقدار یکسان

- به دلیل سرریز شمارنده‌ها، مانند فیلتر بلوم شمارنده

نداشتن مشکل در فیلتر فاخته غیرقابل حذف

- فضاکارآمد

- عدم نیاز به به خاطر سپردن مقادیر یکسانی که چندین بار اضافه شده‌اند

فیلتر فاخته

مثبت کاذب امکان پذیر

مقادیر مختلف دارای اثر انگشت یکسان

- در بیشتر موارد سطل‌های نامزد متفاوتی
- در نتیجه امکان تمایز آنها

با این حال، سطل‌های نامزد یکسان برای آن مقادیر

- رخ دادن تصادم سخت

امکان پاسخ مثبت کاذب به دلیل چنین رویدادهای بسیار نادری

- احتمال مثبت کاذب

$$pr_{fp} = 1 - \left(1 - \frac{1}{2^p}\right)^{2b} \approx \frac{2b}{2^p}$$

فیلتر فاخته

با تعداد ثابت عناصر مورد انتظار n ،

- سبک سنگینی بین احتمال مثبت کاذب pr_{fp} و اندازه سطل b
- جبران با طول اثر انگشتها p
- اثر انگشتها به اندازه کافی طولانی باشند، درهم‌ساز فاخته با کلید جزئی تقریب خوبی از درهم‌ساز فاخته استاندارد
- تاثیر اثر انگشتهای طولانی‌تر بر فضای مورد نیاز
- برآورد طول اثر انگشت توصیه شده p

$$p \geq \left\lceil \log \frac{2b}{pr_{fp}} \right\rceil$$

فیلتر فاخته

کران پائین ذخیره مقادیر به اندازه تعداد مقادیر ورودی در m سطل با اندازه b :

$$m \geq \left\lceil \frac{n}{b} \right\rceil$$

فیلتر فاخته

منفی کاذب امکان ناپذیر

مانند سایر موارد،

- در صورت برگشت عدم عضویت
- قطعاً عدم عضویت

$$pr_{fn} = 0$$

تضمین اشغال فضای بالا

- اصلاح تصمیم‌گیری‌های قبلی قرارگیری مقادیر با درج مقدار جدید
- با این حال، دارای حداکثر ظرفیتی بیان شده با ضریب α
- پس از رسیدن به حداکثر ضریب بار قابل اجرا،
- شکست فزاینده درج‌ها
- نیاز به بزرگ کردن جدول درهم جهت ذخیره مقادیر بیشتر

فیلتر فاخته

- میانگین تعداد بیت‌ها به ازای هر عنصر β
- به عنوان نسبت بین طول اثر انگشت‌ها و ضریب بار α
- تخمین آن با احتمال مثبت کاذب ثابت pr_{fp}

$$\beta \leq \frac{1}{\alpha} \times \left\lceil \log \frac{2b}{pr_{fp}} \right\rceil$$

- استفاده درهم‌ساز فاخته از دو تابع درهم‌ساز
- ضریب بار با سطل‌هایی با اندازه $b = 1$ برابر با پنجاه درصد
- افزایش اندازه سطل امکان بهبود اشغال جدول
- مثال، برای $b = 2$ و $b = 4$ ضریب بار به ترتیب ۸۴ درصد و ۹۵ درصد است.

مطالعه تجربی

- کفایت سطل‌هایی با اندازه $b \in \{1,2,3,4\}$

فیلتر فاخته

مثال - حافظه مورد نیاز

- مدیریت یک میلیارد مقدار با فیلتر فاخته
- احتمال رویدادهای مثبت کاذب در حدود دو درصد
- با اشغال ۸۴ درصدی جدول را ۸۴

نیازمندی‌های پشتیبانی از چنین ضریب باری

- انتخاب اندازه سطل‌ها $b = 2$
- پس طول فیلتر $m = 2^{29}$
- حداقل طول اثر انگشت

$$p = \left\lceil \log \frac{4}{0.02} \right\rceil = 8$$

فیلتر فاخته

بنابراین، طول مورد نیاز اثر انگشت ۸ بیت است

اندازه کل فیلتر فاخته $2^{29} \times 8 \times 2$ بیت

▪ تقریباً معادل ۱,۰۷ گیگابایت حافظه

برای مقایسه نیازهای فضایی با سایر فیلترهای مطالعه شده،

▪ امکان استفاده از $b = 1$

▪ دستیابی به ۵۰ درصد اشغال جدول

▪ نیاز به فیلتری به طول $m = 2^{30}$

▪ استفاده از اثر انگشت‌های ۹ بیتی استفاده کنیم

منجر به حدود ۰,۹۴ گیگابایت حافظه

فیلتر فاخته

استفاده فیلترهای فاخته از رویکردی مشابه فیلتر بلوم شمارنده
▪ با کارایی فضایی بهتری و پیاده‌سازی بسیار ساده‌تر

مناسب برنامه‌های با مقادیر ذخیره شده زیاد و میزان مثبت کاذب نسبتاً پایینی (کمتر از ۳ درصد)
فضای کمتر نسبت به فیلترهای بلوم

زمانی که فیلتر فاخته در حداکثر ظرفیت خود

▪ نیاتز به بزرگ کردن جدول درهم

▪ عدم امکان درج عضو جدید تا آن زمان امکان

▪ در مقابل، درج در فیلتر بلوم به قیمت افزایش نسبت مثبت کاذب