

الگوریتم پردازش استریم DHS آیتم‌های ورودی را به باکت‌ها هاش می‌کند و خانه^۱ (سلول‌ها یا همان شمارنده‌ها) درون هر باکت (اسلات حافظه) را به صورت پویا مدیریت می‌کند. در واقع سائز تمامی خانه‌ها درون یک باکت به صورت سازگار با سائز و توزیع جریان واقعی تنظیم شده‌است. بدین صورت که جریان هر چه بزرگتر شد بیت‌های بیشتر و جریان‌های کوچک بیت‌های کمتر تخصیص می‌یابد. در ابتدا نیز تعداد بیت‌ها کم می‌باشد (اذا مبتنی بر اسکچ می‌باشد و از ویژگی نرخ گذر بالای آن استفاده می‌کند). همچنین یک مکانیزم کوثری **longest fingerprint first** برای تسریع بازیابی استفاده می‌کند.

در بحث ارزیابی، دقت بالا (روش‌های قبلی تخمینی می‌باشند)، نرخ گذر بالا (این که داده‌ها را با سرعتی حداقل برابر با سرعت تولید آنها پردازش کند)، عمومیت بالا و کارایی حافظه DHS را در مقایسه با دیگر راه‌حل‌های هیبریدی ارائه شده، برای مسائل مختلف دیتا استریمینگ (اندازه‌گیری) اثبات می‌کند (تمامی الگوریتم‌ها به زبان C پیاده‌سازی شده و در گیت‌هاب موجود می‌باشند):

- تخمین سائز جریان – k-top
- تخمین توزیع اندازه جریان: تعداد جریان‌هایی که یک سائز مشخص را دارند.
- تشخیص heavy hitter ها
- تشخیص آنومالی‌ها و تغییرات شاخص
- تشخیص آنروپی: آنروپی اندازه جریان

الگوریتم‌های محاسبات دیتا استریمینگ:

با توجه به خصوصیات دیتا استریم‌ها، داده را در حافظه‌های محدود به صورت موقت ذخیره می‌کنند، آنها را یک بار پردازش کرده و آمارش را در داده ساختارهایی برای یک کوپری که آخر یک بازه یا در انتهای کار الگوریتم اجرا می‌شود، نگهداری می‌کنند. داده‌ساختارهای مورد استفاده در دیتا استریمینگ، بایستی یک تخمین کلی در مورد جریان‌ها ارائه دهند. هر آیتم ورودی استریم متعلق به یک جریان می‌باشد. در محاسبات دیتا استریمی با دو متد اصلی **insert** و **query** مواجه هستیم و برای مسائل استخراج داده^۲ پیچیده‌تر، می‌توان توابعی افزون بر اینها و بر پایه اینها ایجاد کرد

هر کدام از الگوریتم‌های دیتا استریمینگ برای یک دسته از محاسبات مناسب (۵ دسته بالایی) می‌باشند و کارایی را در زمینه‌های کارایی حافظه، دقت یا نرخ گذر بهبود می‌بخشند. این روش‌ها را به سه دسته می‌توان تقسیم کرد:

- روش‌های مبتنی بر اسکچ: اطلاعاتی کلی از همه جریان‌ها را ذخیره می‌کنند و از نرخ گذر بالایی برخوردار می‌باشند.
- روش‌های مبتنی بر شمارنده: اطلاعات را به صورت دقیق ذخیره کرده و دارای نرخ گذر پایین می‌باشند. آرایه‌ای از چندین اسلات که هر اسلات یک دوتایی (کلید، مقدار) (یا همان شمارنده) می‌باشد. کلید همان شناسه جریان و مقدار، فرکانس جریان می‌باشد. زمانی که یک جریان جدید می‌آید، جریان با کمترین فرکانس از آرایه خارج خواهد شد.

¹ cells

² Data minning

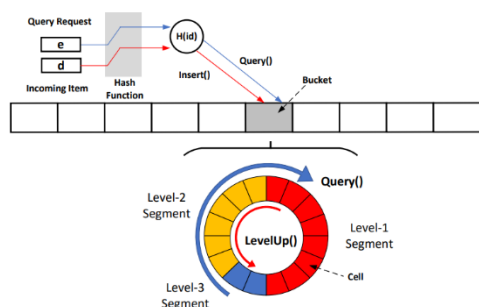
- روش‌های هیبریدی: ترکیبی از ایده‌های اصلی روش‌های بالا با یک مصالحه^۳ بین دقت و نرخ‌گذر(هدف برخورداری از هردو مزیت می‌باشد) می‌باشند و برای مسائل بیشتر اندازه‌گیری مناسب می‌باشند. وبه جای مقدار فرکانس که در اسکچ‌ها استفاده می‌شد، از دوتایی (کلید،مقدار) استفاده می‌کنند. آرایه‌ای از باکت‌ها می‌باشند که هر باکت شامل چندین جفت دودویی است. هر جریان به یک باکت هش می‌شود و (شناسه جریان، فرکانس جریان) در یک باکت ذخیره خواهد شد.

متود ما سه ایده کلی دارد:

- تخصیص سازگار حافظه برای شمارش فرکانس: در شمارنده‌ها و اسکچ‌ها برای جریان‌های ضعیف هم از مقدار حافظه معین تعریف شده، برای اسلات‌ها و ذخیره داده‌ها استفاده می‌کنند. اما در DHS ابتدا سایز کم می‌دهد و سپس با ورود جریان‌های جدید، می‌تواند افزایش یابد.
- تخصیص سازگار حافظه برای شناسه جریان: اسکچ‌ها شناسه جریان را نگهداری نمی‌کنند و منجر به خطای تخمین برای جریان‌های موشی می‌شوند. در طول پردازش آنلاین، زمانی که فرکانس یک جریان پایین است، اثر انگشت^۴ آن با بیت‌های کمتری تخصیص داده می‌شود. همانطور که فرکانس آن افزایش می‌یابد، اثر انگشت آن در محدوده بزرگتری دوباره محاسبه می‌شود.
- تخصیص سازگار بخش‌های از حافظه به جریان‌های موشی و فیلی: در ابتدا همه جریان‌ها موشی اند اما با گذشت زمان برای جریان‌های فیلی نیاز به حافظه بیشتر برای ذخیره سازی شناسه و مقادیر فرکانس جریان داریم که این‌ها را می‌توانیم از موش‌ها بگیریم.

معماری و ساختار:

داده ساختار DHS: داده ساختار ما آرایه ای شامل B باکت که هر باکت شامل یک آرایه مدور^۵ و متدیتا هست، می‌باشد. آرایه مدور به چندین



سگمنت متوالی و سلسله مراتبی منطقی تقسیم شده‌است. هر سگمنت چندین خانه (سلول) پایه‌ای دارد که یک دوتایی (اثر انگشت، فرکانس) می‌باشد اما سایز این دوتایی بنا به سگمنت آن می‌تواند متغیر باشد. مرز بین سگمنت‌ها قابل تغییر می‌باشد و در متادیتای باکت ذخیره شده‌است. هر باکت یک آرایه ۳ سطحی مدور می‌باشد. فیلدهای اثر انگشت و فرکانس، تعداد بیت‌هایشان یکسان می‌باشد و برای سطوح ۱، ۲ و ۳ به ترتیب ۸، ۱۲ و ۱۶ بیت می‌باشد. طول آرایه مدور (سایز باکت) عدد ثابت W می‌باشد. تابع هش h0 برای محاسبه اندیس جریان (0, B-1) در آرایه باکت می‌باشد. h1, h2, h3 اثر انگشت هر سطح را برمی‌گردانند و بازه‌های

آنها به ترتیب (0-255)، (0-4095) و (0-65535) می‌باشد. در حین پردازش استریم، مرز هر سگمنت بنا به جریان داخلش تغییر می‌یابد، که در نهایت منجر به تنظیم خانه‌ها (سلول‌ها) در هر سطح می‌شود. باید توجه داشت که زمانی که یک باکت سرشار از خانه‌های سطح بالا باشد، هر باکت مثل روش شمارنده عمل می‌کند و اگر سرشار از خانه‌های سطح پایین باشد، همانند یک اسکچ عمل خواهد کرد.

طریقه کار: در ابتدا همه فیلدها برابر صفر و هر باکت تنها خانه‌های سطح پایین می‌باشد. شامل عملیات زیر می‌باشد:

- query: از روش longest fingerprint first استفاده می‌کند. ابتدا باکت متناظرش و سپس خانه‌ها را از سطح پایین تا سطح بالا پیمایش می‌کند. خروجی توابع هش را بر اساس شناسه جریان محاسبه می‌کند.
- insert: پس از آن که باکت متناظر را پیدا کرد، در آن باکت عمل query را اجرا می‌کند که حالت‌های زیر ممکن است رخ دهد:

³ Trade-off

⁴ fingerprint

⁵ Circular array

⁶ cell

- اگر دریکی از سطحها رکورد پیدا شد و فرکانس مورد نظر را افزایش میدهد با این که هیچ گونه سرریزی رخ نخواهد داد.
- اگر دریکی از سطحها (k) رکورد پیدا شد اما در صورت درج، سرریز رخ می‌داد، از سطح 0 تا k برای یافتن سطحی که خانه آزاد کافی برای سطح جدید k+1 را داشته باشد، جستجو می‌کند. اگر جستجو موفق بود، خانه‌های آزاد شده را به سطح k+1 منتقل می‌کند. اگر موفق نبود، با یک احتمالی^۷ از مقدار فرکانس کوچکترین خانه سطح k+1 می‌کاهد و اگر مقدار آن از فرکانس جریان ورودی کمتر بود، با جریان ورودی جایگزین خواهد شد.
- اگر query پاسخ منفی برگرداند، خانه جدید به سطح یک اضافه خواهد شد.
- اگر query پاسخ منفی برگرداند و سطح یک خانه جدید نداشت، با یک احتمالی^۷ از فرکانس کوچکترین خانه سطح ۱ می‌کاهد و اگر مقدار کوچکترین خانه به ۰ رسید، جریان ورودی را با آن جایگذاری می‌کند.

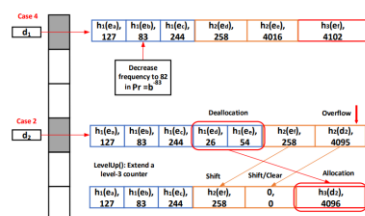


Figure 3: Insertion Examples of DHS

ارزیابی راهکار ارائه شده و مقایسه با دیگر راهکار:

پیش از مقایسه DHS با دیگر راهکارها، به بررسی میزان تاثیر پارامترهای تعداد باکت (B)، اندازه باکت (W) و پارامتر exponential decay (b) پرداخته و بهینه ترین آن‌ها را انتخاب می‌کند.

⁷ Exponential decay parameter