

PHast

Perfect Hashing made fast

Piotr Beling i Peter Sanders

2026

Na podstawie artykułu

Piotr Beling, Peter Sanders *PHast - Perfect Hashing made fast* 2026 Proceedings of the SIAM Symposium on Algorithm Engineering and Experiments (ALENEX), którego pełna, najnowsza wersja dostępna jest na <https://arxiv.org/abs/2504.17918>, oraz slajdów pokazanych na konferencji ALENEX 2026.

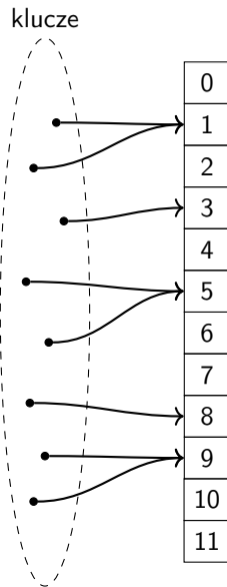
Perfect hashing

Dla zadanego zbioru kluczy $\{k_1, \dots, k_n\}$, funkcja haszująca (*Hash Function*) $f : \{k_1, \dots, k_n\} \rightarrow \{0, \dots, m - 1\}$ jest

- *Doskonała (Perfect; PHF)*
jeśli jest różnowartościowa,
- *Minimalna Doskonała (Minimal Perfect; MPHf)*
jesli, dodatkowo, $m = n$.

Pożądane cechy MPHf:

- **mały rozmiar**
(teoretyczne minimum to $\log_2 e \approx 1.44$ bitów na klucz,
niezależnie od typów kluczy),
- **maksymalna szybkość**
konstruowania i obliczania wartości.



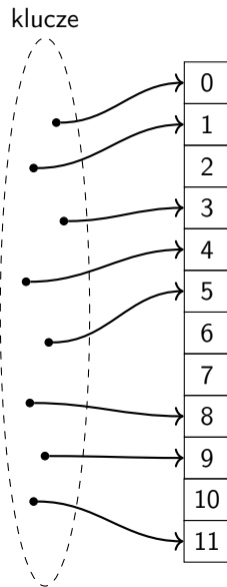
Perfect hashing

Dla zadanego zbioru kluczy $\{k_1, \dots, k_n\}$, funkcja haszująca (*Hash Function*) $f : \{k_1, \dots, k_n\} \rightarrow \{0, \dots, m - 1\}$ jest

- *Doskonała (Perfect; PHF)*
jeśli jest różnowartościowa,
- *Minimalna Doskonała (Minimal Perfect; MPHf)*
jesli, dodatkowo, $m = n$.

Pożądane cechy MPHf:

- **mały rozmiar**
(teoretyczne minimum to $\log_2 e \approx 1.44$ bitów na klucz,
niezależnie od typów kluczy),
- **maksymalna szybkość**
konstruowania i obliczania wartości.



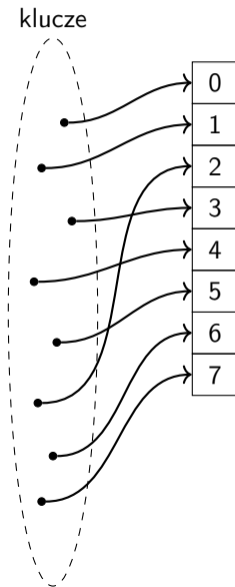
Perfect hashing

Dla zadanego zbioru kluczy $\{k_1, \dots, k_n\}$, funkcja haszująca (*Hash Function*) $f : \{k_1, \dots, k_n\} \rightarrow \{0, \dots, m - 1\}$ jest

- *Doskonała (Perfect; PHF)*
jeśli jest różnowartościowa,
- *Minimalna Doskonała (Minimal Perfect; MPHf)*
jesli, dodatkowo, $m = n$.

Pożądane cechy MPHf:

- **mały rozmiar**
(teoretyczne minimum to $\log_2 e \approx 1.44$ bitów na klucz,
niezależnie od typów kluczy),
- **maksymalna szybkość**
konstruowania i obliczania wartości.



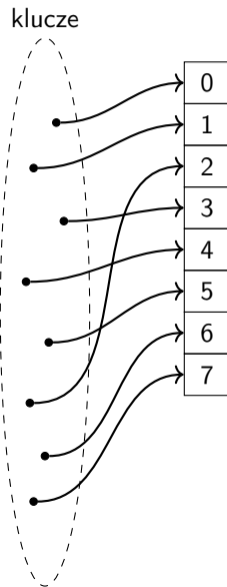
Perfect hashing

Dla zadanego zbioru kluczy $\{k_1, \dots, k_n\}$, funkcja haszująca (*Hash Function*) $f : \{k_1, \dots, k_n\} \rightarrow \{0, \dots, m - 1\}$ jest

- *Doskonała (Perfect; PHF)*
jeśli jest różnowartościowa,
- *Minimalna Doskonała (Minimal Perfect; MPHf)*
jesli, dodatkowo, $m = n$.

Pożądane cechy MPHf:

- **mały rozmiar**
(teoretyczne minimum to $\log_2 e \approx 1.44$ bitów na klucz,
niezależnie od typów kluczy),
- **maksymalna szybkość**
konstruowania i obliczania wartości.



Zastosowania doskonałych funkcji haszujących

(Minimalne) Doskonałe Funkcje Haszujące są używane w:

- bazach danych,
 - indeksowaniu tekstu,
 - statycznych tablicach haszujących,
 - klasyfikacji pakietów w routerach i firewallach,
 - bioinformatyce.
-
- Szybka ewaluacja jest kluczowa w wielu zastosowaniach, i dlatego była głównym celem projektowym PHast (Perfect Hashing made fast).
 - Dlatego PHast bazuje na technice zwanej *umieszczaniem wiader* (*bucket placement*), wykorzystywanej przez najszybsze, znane doskonałe funkcje haszujące.

(Minimalne) Doskonałe Funkcje Haszujące są używane w:

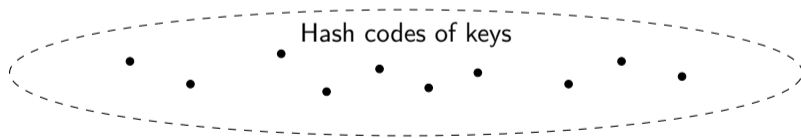
- bazach danych,
 - indeksowaniu tekstu,
 - statycznych tablicach haszujących,
 - klasyfikacji pakietów w routerach i firewallach,
 - bioinformatyce.
-
- **Szybka ewaluacja jest kluczowa** w wielu zastosowaniach, i dlatego była głównym celem projektowym PHast (Perfect Hashing made fast).
 - Dlatego PHast bazuje na technice zwanej *umieszczaniem wiader* (*bucket placement*), wykorzystywanej przez najszybsze, znane doskonałe funkcje haszujące.

(Minimalne) Doskonałe Funkcje Haszujące są używane w:

- bazach danych,
 - indeksowaniu tekstu,
 - statycznych tablicach haszujących,
 - klasyfikacji pakietów w routerach i firewallach,
 - bioinformatyce.
-
- **Szybka ewaluacja jest kluczowa** w wielu zastosowaniach, i dlatego była głównym celem projektowym PHast (Perfect Hashing made fast).
 - Dlatego PHast bazuje na technice zwanej *umieszczaniem wiader* (*bucket placement*), wykorzystywanej przez najszybsze, znane doskonałe funkcje haszujące.

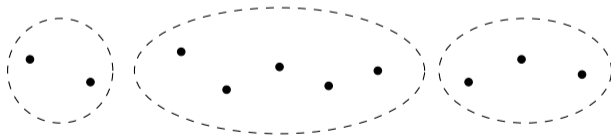
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

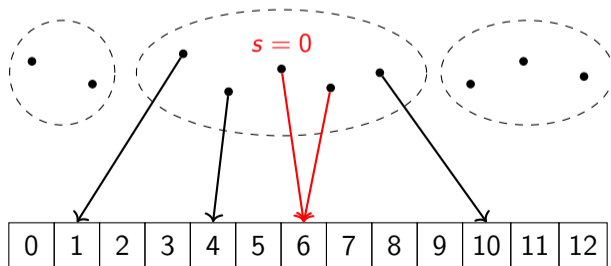
- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



0	1	2	3	4	5	6	7	8	9	10	11	12
---	---	---	---	---	---	---	---	---	---	----	----	----

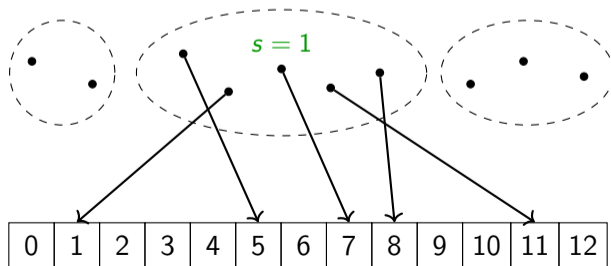
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



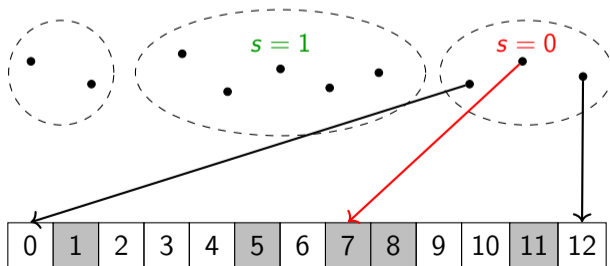
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



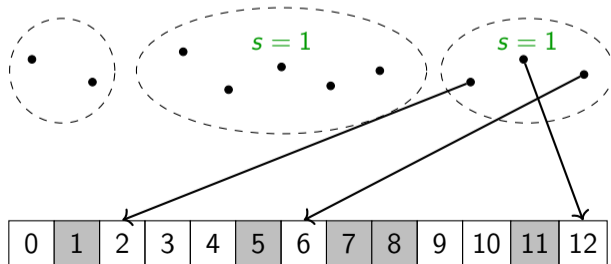
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



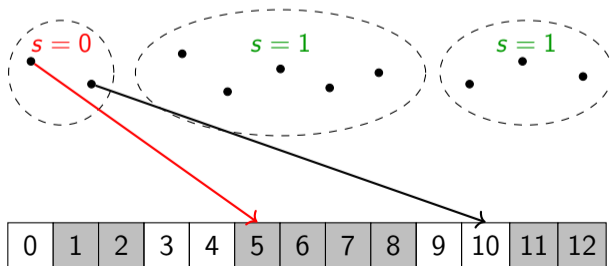
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



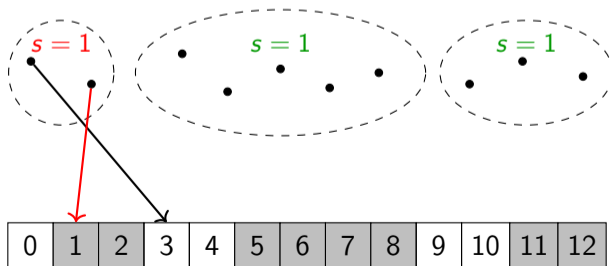
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



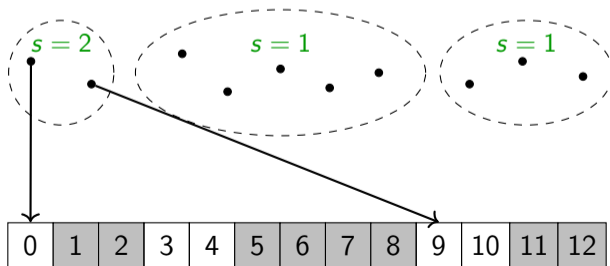
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



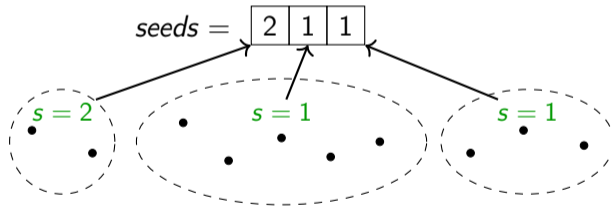
Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



Konstruowanie typowego PHF bezującego na umieszczaniu kubełków

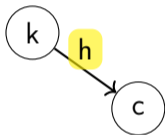
- 1 przypisanie wartości haszującej $c = h(k)$ do każdego klucza k
- 2 pseudolosowy podział zbioru wartości haszujących na kubełki
- 3 dla każdego kubełka, w kolejności od największego do najmniejszego:
znalezienie najmniejszego załączka s dla którego funkcja $p(s, c)$ bezkolizyjnie przypisuje wartości do wszystkich c w kubełku
- 4 zapis znalezionych załączków w (opcjonalnie skompresowanej) tablicy *seeds*



0	1	2	3	4	5	6	7	8	9	10	11	12
---	---	---	---	---	---	---	---	---	---	----	----	----

Obliczanie (ewaluacja) $f(k)$:

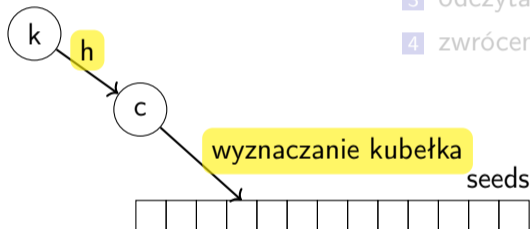
- 1 obliczenie wartości haszującej $c = h(k)$;
- 2 wyznaczenie indeksu i kubełka zawierającego c ;
- 3 odczytanie załączka $s = seeds[i]$ i -tego kubełka;
- 4 zwrócenie $p(s, c)$.



Ewaluacja typowego PHF bazującego na umieszczaniu kubełków

Obliczanie (ewaluacja) $f(k)$:

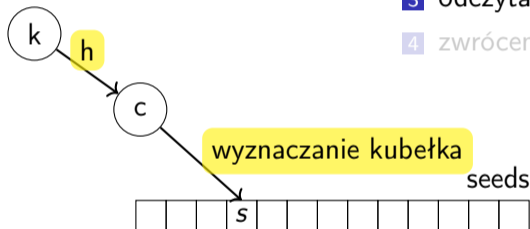
- 1 obliczenie wartości haszującej $c = h(k)$;
- 2 wyznaczenie indeksu i kubełka zawierającego c ;
- 3 odczytanie załączka $s = seeds[i]$ i -tego kubełka;
- 4 zwrócenie $p(s, c)$.



Ewaluacja typowego PHF bazującego na umieszczaniu kubełków

Obliczanie (ewaluacja) $f(k)$:

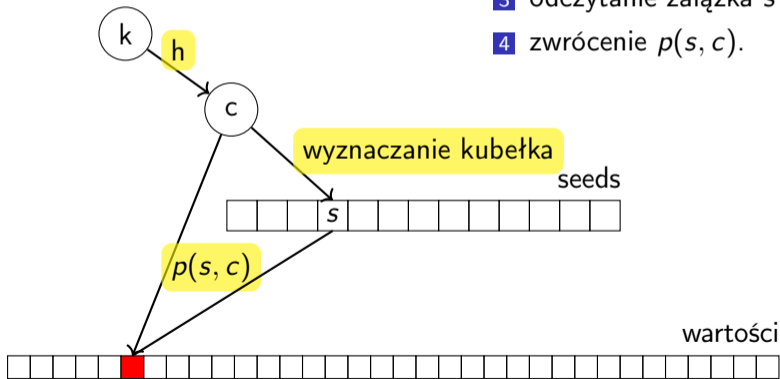
- 1 obliczenie wartości haszującej $c = h(k)$;
- 2 wyznaczenie indeksu i kubełka zawierającego c ;
- 3 odczytanie załączka $s = seeds[i]$ i -tego kubełka;
- 4 zwrócenie $p(s, c)$.



Ewaluacja typowego PHF bazującego na umieszczaniu kubełków

Obliczanie (ewaluacja) $f(k)$:

- 1 obliczenie wartości haszującej $c = h(k)$;
- 2 wyznaczenie indeksu i kubełka zawierającego c ;
- 3 odczytanie załączka $s = seeds[i]$ i -tego kubełka;
- 4 zwrócenie $p(s, c)$.



Załączki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- Kodowanie o zmiennej długości jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Załączki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- **Kodowanie o zmiennej długości** jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Zalążki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- **Kodowanie o zmiennej długości** jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Zalążki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- **Kodowanie o zmiennej długości** jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Załączki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- **Kodowanie o zmiennej długości** jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Załączki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- **Kodowanie o zmiennej długości** jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Zalążki dla ostatnich kubełków są typowo znacznie wyższe bo trudniej je znaleźć.

- **Kodowanie o zmiennej długości** jest niezbędne do uzyskania małego rozmiaru.
- By wyrównać trudność znalezienia załączków, PHFy takie jak PHOBIC i PTHash mocniej różnicują wielkości kubełków używając **nieliniowego przypisywania kubełków** do wartości haszujących.

W razie kolizji wartości haszujących ($h(k_1) = h(k_2)$ dla $k_1 \neq k_2$), konieczne jest **ponowne haszowanie** kluczy za pomocą innego h .

- Ze względu na paradoks dnia urodzin, h musi mieć duży zbiór wartości.

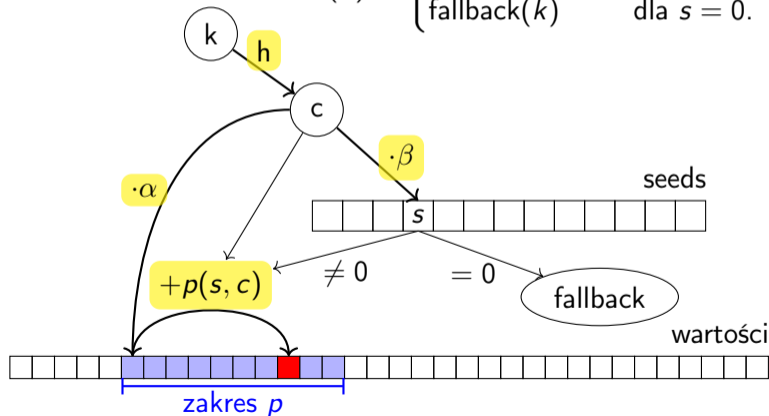
Konstruowanie wielowątkowe wymaga **podziału zbioru wejściowego na partycje** (po jednej na wątek), co wiąże się z koniecznością odczytu i dodania przesunięcia partycji podczas ewaluacji.

Wszystko to spowalnia ewaluację.

Obliczanie wartości przez PHast

Wartość PHast dla klucza k o haszu $c = h(k)$ i załączka $s = seeds[i]$ (gdzie $i = \lfloor \beta c \rfloor$ jest indeksem kubetka) to

$$f(k) = \begin{cases} \lfloor \alpha c \rfloor + p(s, c) & \text{dla } s \neq 0, \\ \text{fallback}(k) & \text{dla } s = 0. \end{cases}$$

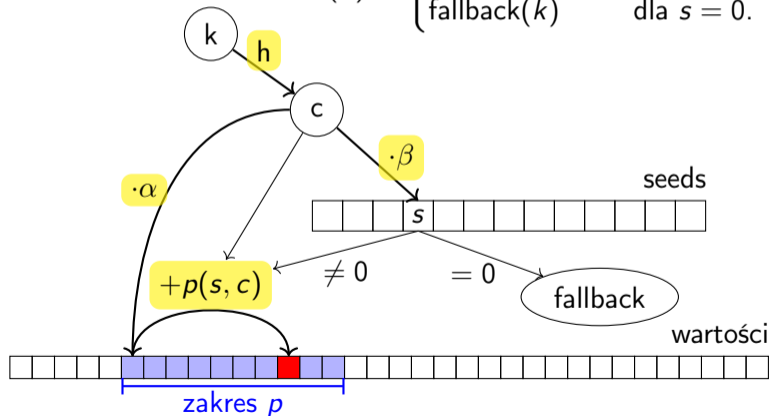


- $s = 0$ oznacza **wybicie** kubetka (dla którego nie znaleziono załączka)
- dzięki wybijaniu mamy:
 - stały rozmiar załączków
 - h o małym zakresie
 - proste (i szybkie) p
- dzięki małemu zakresowi p , i -ty kubetek pokrywa ograniczony zakres, przesuwany ku większym wartościom gdy wzrasta i
- to czyni konstruowanie przyjaznym dla CPU-cache i wielowątkowości

Obliczanie wartości przez PHast

Wartość PHast dla klucza k o haszu $c = h(k)$ i załączka $s = seeds[i]$ (gdzie $i = \lfloor \beta c \rfloor$ jest indeksem kubetka) to

$$f(k) = \begin{cases} \lfloor \alpha c \rfloor + p(s, c) & \text{dla } s \neq 0, \\ \text{fallback}(k) & \text{dla } s = 0. \end{cases}$$

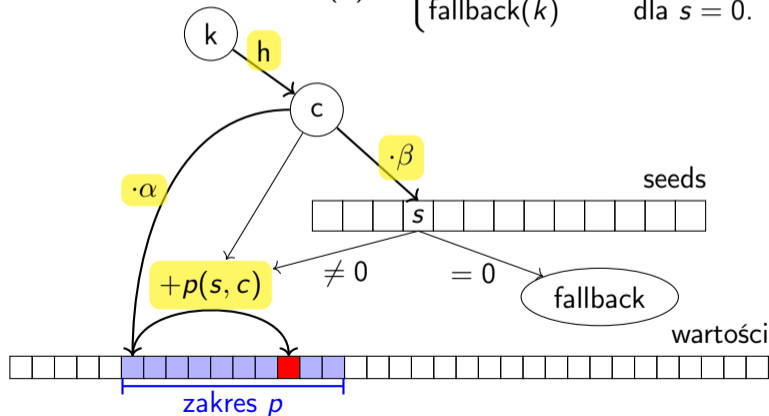


- $s = 0$ oznacza **wybicie** kubetka (dla którego nie znaleziono załączka)
- dzięki wybijaniu mamy:
 - stały rozmiar załączków
 - h o małym zakresie
 - proste (i szybkie) p
- dzięki małemu zakresowi p , i -ty kubetek pokrywa ograniczony zakres, przesuwany ku większym wartościom gdy wzrasta i
- to czyni konstruowanie przyjaznym dla CPU-cache i wielowątkowości

Obliczanie wartości przez PHast

Wartość PHast dla klucza k o haszu $c = h(k)$ i załączka $s = seeds[i]$ (gdzie $i = \lfloor \beta c \rfloor$ jest indeksem kubetka) to

$$f(k) = \begin{cases} \lfloor \alpha c \rfloor + p(s, c) & \text{dla } s \neq 0, \\ \text{fallback}(k) & \text{dla } s = 0. \end{cases}$$

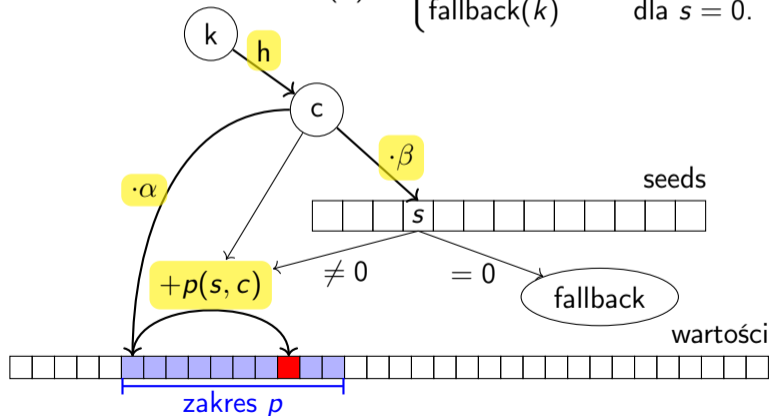


- $s = 0$ oznacza **wybicie** kubetka (dla którego nie znaleziono załączka)
- dzięki wybijaniu mamy:
 - stały rozmiar załączków
 - h o małym zakresie
 - proste (i szybkie) p
- dzięki małemu zakresowi p , i -ty kubetek pokrywa ograniczony zakres, przesuwany ku większym wartościom gdy wzrasta i
- to czyni konstruowanie przyjaznym dla CPU-cache i wielowątkowości

Obliczanie wartości przez PHast

Wartość PHast dla klucza k o haszu $c = h(k)$ i załączka $s = seeds[i]$ (gdzie $i = \lfloor \beta c \rfloor$ jest indeksem kubetka) to

$$f(k) = \begin{cases} \lfloor \alpha c \rfloor + p(s, c) & \text{dla } s \neq 0, \\ \text{fallback}(k) & \text{dla } s = 0. \end{cases}$$

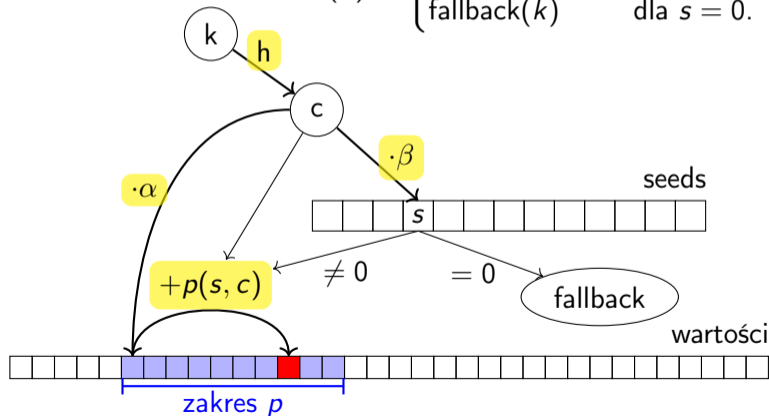


- $s = 0$ oznacza **wybicie** kubetka (dla którego nie znaleziono załączka)
- dzięki wybijaniu mamy:
 - stały rozmiar załączków
 - h o małym zakresie
 - proste (i szybkie) p
- dzięki małemu zakresowi p , i -ty kubetek pokrywa ograniczony zakres, przesuwany ku większym wartościom gdy wzrasta i
- to czyni konstruowanie przyjaznym dla CPU-cache i wielowątkowości

Obliczanie wartości przez PHast

Wartość PHast dla klucza k o haszu $c = h(k)$ i załączka $s = seeds[i]$ (gdzie $i = \lfloor \beta c \rfloor$ jest indeksem kubetka) to

$$f(k) = \begin{cases} \lfloor \alpha c \rfloor + p(s, c) & \text{dla } s \neq 0, \\ \text{fallback}(k) & \text{dla } s = 0. \end{cases}$$



- $s = 0$ oznacza **wybicie** kubetka (dla którego nie znaleziono załączka)
- dzięki wybijaniu mamy:
 - stały rozmiar załączków
 - h o małym zakresie
 - proste (i szybkie) p
- dzięki małemu zakresowi p , i -ty kubetek pokrywa ograniczony zakres, przesuwany ku większym wartościom gdy wzrasta i
- to czyni konstruowanie przyjaznym dla CPU-cache i wielowątkowości

załążki dla kubeków:

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

- Kubki są przetwarzane (przypisywane są do nich załączki) mniej więcej po kolei, zgodnie z kolejnością ich indeksów, od lewej do prawej.
- W konsekwencji, wartości funkcji też są zajmowane mniej więcej od lewej do prawej.

załączki dla kubeków:

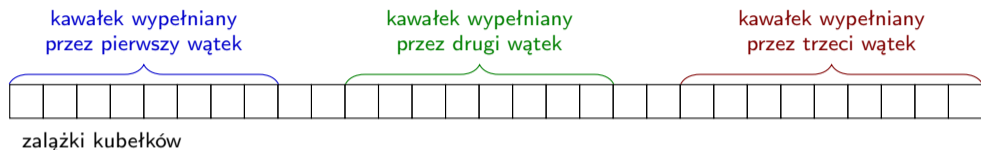
4	0	3	6	2	7	5	1		5		2								
---	---	---	---	---	---	---	---	--	---	--	---	--	--	--	--	--	--	--	--

zajęte/wolne wartości:



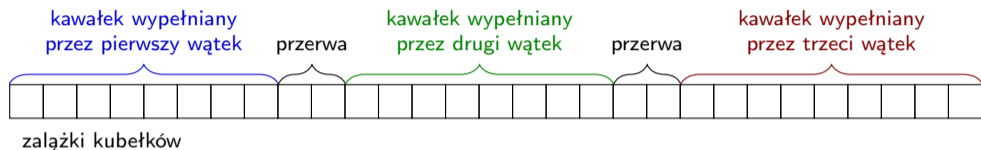
- Kubki są przetwarzane (przypisywane są do nich załączki) mniej więcej po kolei, zgodnie z kolejnością ich indeksów, od lewej do prawej.
- W konsekwencji, wartości funkcji też są zajmowane mniej więcej od lewej do prawej.

Konstruowanie wielowątkowe



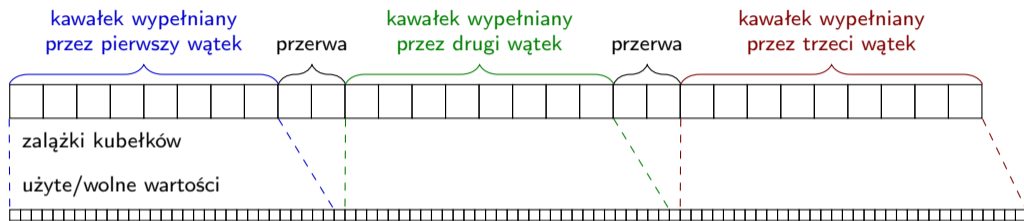
- tablica załączków jest podzielona na kawałki wypełniane przez osobne wątki,
- i (niewielkie) przerwy pomiędzy tymi kawałkami, wypełniane później;
- dzięki tym przerwom:
- wartości funkcji pokrywane przez kawałki nie zazębiają się,
- dzięki czemu nie jest potrzebna żadna komunikacja pomiędzy wątkami

Konstruowanie wielowątkowe



- tablica załączków jest podzielona na kawałki wypełniane przez osobne wątki,
- i (niewielkie) przerwy pomiędzy tymi kawałkami, wypełniane później;
- dzięki tym przerwom:
- wartości funkcji pokrywane przez kawałki nie zazębiają się,
- dzięki czemu nie jest potrzebna żadna komunikacja pomiędzy wątkami

Konstruowanie wielowątkowe



- tablica załączków jest podzielona na kawałki wypełniane przez osobne wątki,
- i (niewielkie) przerwy pomiędzy tymi kawałkami, wypełniane później;
- dzięki tym przerwom:
- wartości funkcji pokrywane przez kawałki nie zazębiają się,
- dzięki czemu nie jest potrzebna żadna komunikacja pomiędzy wątkami

Szczegóły implementacyjne

- Nasza implementacja PHasta jest rozwijana w Ruście i może używać dowolnej funkcji 64-bitowej funkcji haszującej h .
- Domyślnie używany jest GxHash, który odznacza się dużą szybkością i jakością.
- Regularny PHast używa następującej funkcji rozmieszczającej

$$p(s, c) = (\lfloor s \cdot 5871781006564002453 / 2^{64} \rfloor \cdot c) \bmod L,$$

gdzie L ogranicza zakres p (by móc szybko wykonywać dzielenie modulo za pomocą *bitowego and*, używamy jedynie L będących potęgami dwójki).

- PHast⁺ – wariant z szybkim (bitowo-zrównoległym) konstruowaniem – używa:

$$p(s, c) = c \bmod L + s - 1, \tag{1}$$

$$p(s, c) = (c + \delta s) \bmod L, \tag{2}$$

gdzie δ jest parametrem metody (używamy $\delta \in \{1, 2, 3\}$).

Szczegóły implementacyjne

- Nasza implementacja PHasta jest rozwijana w Ruście i może używać dowolnej funkcji 64-bitowej funkcji haszującej h .
- Domyślnie używany jest GxHash, który odznacza się dużą szybkością i jakością.
- Regularny PHast używa następującej funkcji rozmieszczającej

$$p(s, c) = (\lfloor s \cdot 5871781006564002453 / 2^{64} \rfloor \cdot c) \bmod L,$$

gdzie L ogranicza zakres p (by móc szybko wykonywać dzielenie modulo za pomocą *bitowego and*, używamy jedynie L będących potęgami dwójki).

- PHast⁺ – wariant z szybkim (bitowo-zrównoległym) konstruowaniem – używa:

$$p(s, c) = c \bmod L + s - 1, \tag{1}$$

$$p(s, c) = (c + \delta s) \bmod L, \tag{2}$$

gdzie δ jest parametrem metody (używamy $\delta \in \{1, 2, 3\}$).

Szczegóły implementacyjne

- Nasza implementacja PHasta jest rozwijana w Ruście i może używać dowolnej funkcji 64-bitowej funkcji haszującej h .
- Domyślnie używany jest GxHash, który odznacza się dużą szybkością i jakością.
- Regularny PHast używa następującej funkcji rozmieszczającej

$$p(s, c) = (\lfloor s \cdot 5871781006564002453/2^{64} \rfloor \cdot c) \bmod L,$$

gdzie L ogranicza zakres p (by móc szybko wykonywać dzielenie modulo za pomocą *bitowego and*, używamy jedynie L będących potęgami dwójki).

- PHast⁺ – wariant z szybkim (bitowo-zrównoległym) konstruowaniem – używa:

$$p(s, c) = c \bmod L + s - 1, \tag{1}$$

$$p(s, c) = (c + \delta s) \bmod L, \tag{2}$$

gdzie δ jest parametrem metody (używamy $\delta \in \{1, 2, 3\}$).

Szczegóły implementacyjne

- Nasza implementacja PHasta jest rozwijana w Ruście i może używać dowolnej funkcji 64-bitowej funkcji haszującej h .
- Domyślnie używany jest GxHash, który odznacza się dużą szybkością i jakością.
- Regularny PHast używa następującej funkcji rozmieszczającej

$$p(s, c) = (\lfloor s \cdot 5871781006564002453 / 2^{64} \rfloor \cdot c) \bmod L,$$

gdzie L ogranicza zakres p (by móc szybko wykonywać dzielenie modulo za pomocą *bitowego and*, używamy jedynie L będących potęgami dwójki).

- PHast⁺ – wariant z szybkim (bitowo-zrównoległym) konstruowaniem – używa:

$$p(s, c) = c \bmod L + s - 1, \tag{1}$$

$$p(s, c) = (c + \delta s) \bmod L, \tag{2}$$

gdzie δ jest parametrem metody (używamy $\delta \in \{1, 2, 3\}$).

Niech C będzie multi-zbiorem wartości haszujących kluczy w kubełku.

Dla $p(s, c) = c \bmod L + s - 1$, załączek s minimalizujący $\sum_{c \in C} p(s, c)$ można szybko znaleźć, gdyż:

- ponieważ wartości p rosną z s ($p(s + 1, c) = p(s, c) + 1$), wystarczy znaleźć najmniejszy bezkolizyjny załączek;
- test kolizji pomiędzy kluczami w kubełku wystarczy wykonać raz, bo jego wynik nie zależy od s ;
- można zrównoleglić sprawdzanie bezkolizyjności załączków za pomocą operacji bitowych.

Dla $p(s, c) = (c + \delta s) \bmod L$:

- powyższe fakty zachodzą i pomagają znaleźć optymalne załączki osobno w przedziałach, takich że:
załączki a, b są w tym samym przedziale $\Leftrightarrow \lfloor \frac{c+\delta a}{L} \rfloor = \lfloor \frac{c+\delta b}{L} \rfloor$ dla wszystkich $c \in C$;
- większe δ to większa liczba przedziałów i, co za tym idzie, czas konstruowania, ale też mniejszy rozmiar skonstruowanej funkcji.

Szybkie konstruowanie PHast⁺

Niech C będzie multi-zbiorem wartości haszujących kluczy w kubełku.

Dla $p(s, c) = c \bmod L + s - 1$, załączek s minimalizujący $\sum_{c \in C} p(s, c)$ można szybko znaleźć, gdyż:

- ponieważ wartości p rosną z s ($p(s + 1, c) = p(s, c) + 1$), wystarczy znaleźć najmniejszy bezkolizyjny załączek;
- test kolizji pomiędzy kluczami w kubełku wystarczy wykonać raz, bo jego wynik nie zależy od s ;
- można zrównoleglić sprawdzanie bezkolizyjności załączków za pomocą operacji bitowych.

Dla $p(s, c) = (c + \delta s) \bmod L$:

- powyższe fakty zachodzą i pomagają znaleźć optymalne załączki osobno w przedziałach, takich że:
załączki a, b są w tym samym przedziale $\Leftrightarrow \lfloor \frac{c+\delta a}{L} \rfloor = \lfloor \frac{c+\delta b}{L} \rfloor$ dla wszystkich $c \in C$;
- większe δ to większa liczba przedziałów i, co za tym idzie, czas konstruowania, ale też mniejszy rozmiar skonstruowanej funkcji.

Szybkie konstruowanie PHast⁺

Niech C będzie multi-zbiorem wartości haszujących kluczy w kubełku.

Dla $p(s, c) = c \bmod L + s - 1$, załączek s minimalizujący $\sum_{c \in C} p(s, c)$ można szybko znaleźć, gdyż:

- ponieważ wartości p rosną z s ($p(s + 1, c) = p(s, c) + 1$), wystarczy znaleźć najmniejszy bezkolizyjny załączek;
- test kolizji pomiędzy kluczami w kubełku wystarczy wykonać raz, bo jego wynik nie zależy od s ;
- można zrównoleglić sprawdzanie bezkolizyjności załączków za pomocą operacji bitowych.

Dla $p(s, c) = (c + \delta s) \bmod L$:

- powyższe fakty zachodzą i pomagają znaleźć optymalne załączki osobno w przedziałach, takich że:
załączki a, b są w tym samym przedziale $\Leftrightarrow \lfloor \frac{c+\delta a}{L} \rfloor = \lfloor \frac{c+\delta b}{L} \rfloor$ dla wszystkich $c \in C$;
- większe δ to większa liczba przedziałów i, co za tym idzie, czas konstruowania, ale też mniejszy rozmiar skonstruowanej funkcji.

Niech C będzie multi-zbiorem wartości haszujących kluczy w kubełku.

Dla $p(s, c) = c \bmod L + s - 1$, załączek s minimalizujący $\sum_{c \in C} p(s, c)$ można szybko znaleźć, gdyż:

- ponieważ wartości p rosną z s ($p(s + 1, c) = p(s, c) + 1$), wystarczy znaleźć najmniejszy bezkolizyjny załączek;
- test kolizji pomiędzy kluczami w kubełku wystarczy wykonać raz, bo jego wynik nie zależy od s ;
- można zrównoleglić sprawdzanie bezkolizyjności załączków za pomocą operacji bitowych.

Dla $p(s, c) = (c + \delta s) \bmod L$:

- powyższe fakty zachodzą i pomagają znaleźć optymalne załączki osobno w przedziałach, takich że:
załączki a, b są w tym samym przedziale $\Leftrightarrow \lfloor \frac{c+\delta a}{L} \rfloor = \lfloor \frac{c+\delta b}{L} \rfloor$ dla wszystkich $c \in C$;
- większe δ to większa liczba przedziałów i, co za tym idzie, czas konstruowania, ale też mniejszy rozmiar skonstruowanej funkcji.

Niech C będzie multi-zbiorem wartości haszujących kluczy w kubełku.

Dla $p(s, c) = c \bmod L + s - 1$, załączek s minimalizujący $\sum_{c \in C} p(s, c)$ można szybko znaleźć, gdyż:

- ponieważ wartości p rosną z s ($p(s + 1, c) = p(s, c) + 1$), wystarczy znaleźć najmniejszy bezkolizyjny załączek;
- test kolizji pomiędzy kluczami w kubełku wystarczy wykonać raz, bo jego wynik nie zależy od s ;
- można zrównoleglić sprawdzanie bezkolizyjności załączków za pomocą operacji bitowych.

Dla $p(s, c) = (c + \delta s) \bmod L$:

- powyższe fakty zachodzą i pomagają znaleźć optymalne załączki osobno w przedziałach, takich że:
załączki a, b są w tym samym przedziale $\Leftrightarrow \lfloor \frac{c+\delta a}{L} \rfloor = \lfloor \frac{c+\delta b}{L} \rfloor$ dla wszystkich $c \in C$;
- większe δ to większa liczba przedziałów i, co za tym idzie, czas konstruowania, ale też mniejszy rozmiar skonstruowanej funkcji.

Niech C będzie multi-zbiorem wartości haszujących kluczy w kubełku.

Dla $p(s, c) = c \bmod L + s - 1$, załączek s minimalizujący $\sum_{c \in C} p(s, c)$ można szybko znaleźć, gdyż:

- ponieważ wartości p rosną z s ($p(s + 1, c) = p(s, c) + 1$), wystarczy znaleźć najmniejszy bezkolizyjny załączek;
- test kolizji pomiędzy kluczami w kubełku wystarczy wykonać raz, bo jego wynik nie zależy od s ;
- można zrównoleglić sprawdzanie bezkolizyjności załączków za pomocą operacji bitowych.

Dla $p(s, c) = (c + \delta s) \bmod L$:

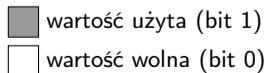
- powyższe fakty zachodzą i pomagają znaleźć optymalne załączki osobno w przedziałach, takich że:
załączki a, b są w tym samym przedziale $\Leftrightarrow \lfloor \frac{c+\delta a}{L} \rfloor = \lfloor \frac{c+\delta b}{L} \rfloor$ dla wszystkich $c \in C$;
- większe δ to większa liczba przedziałów i, co za tym idzie, czas konstruowania, ale też mniejszy rozmiar skonstruowanej funkcji.

Bitowo-zrównległone sprawdzanie załączków w PHast⁺



Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków $(s, s + 1, \dots, s + k - 1)$:

- Dla każdego c w kubelku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.

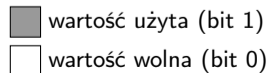


Bitowo-zrównległone sprawdzanie załączków w PHast⁺



Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków $(s, s + 1, \dots, s + k - 1)$:

- Dla każdego c w kubelku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.

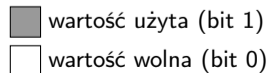
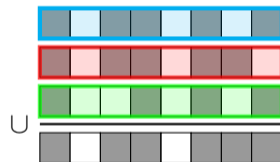


Bitowo-zrównległone sprawdzanie załączków w PHast⁺



Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków $(s, s + 1, \dots, s + k - 1)$:

- Dla każdego c w kubelku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.

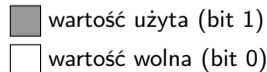
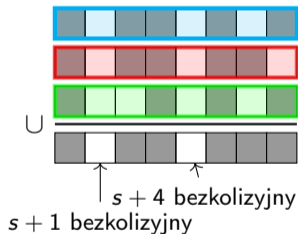


Bitowo-zrównległone sprawdzanie załączków w PHast⁺



Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków ($s, s + 1, \dots, s + k - 1$):

- Dla każdego c w kubeczku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.

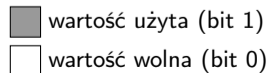
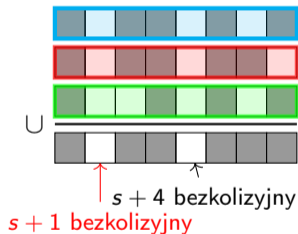


Bitowo-zrównległone sprawdzanie załączków w PHast⁺



Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków ($s, s + 1, \dots, s + k - 1$):

- Dla każdego c w kubeczku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.

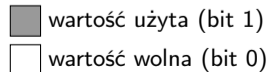
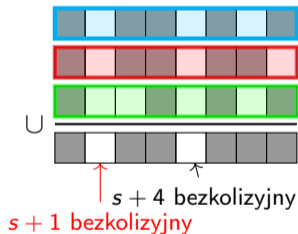


Bitowo-zrównległone sprawdzanie załączków w PHast⁺



Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków ($s, s + 1, \dots, s + k - 1$):

- Dla każdego c w kubeczku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.

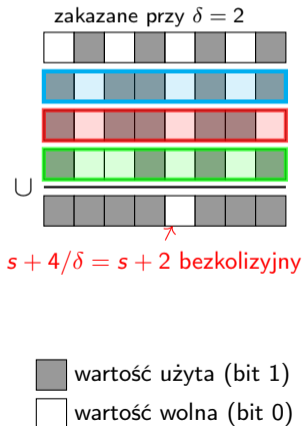


Bitowo-zrównległone sprawdzanie załączków w PHast⁺

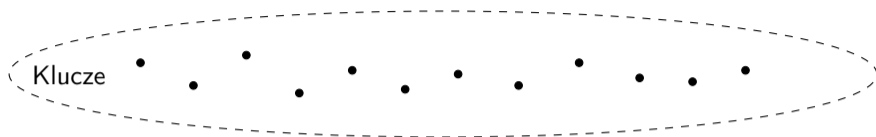


Równoczesne badanie kolizyjności k ($k = 64$ dla współczesnych CPU; $k = 8$ na rysunku) kolejnych załączków ($s, s + 1, \dots, s + k - 1$):

- Dla każdego c w kubeczku, z bitmapy użytych wartości, czytamy k kolejnych bitów poczynając wartości dla pary (s, c) .
- Niech r będzie *bitowym LUB* odczytanych słów.
- i -ty bit r to 0 \Leftrightarrow załączek $s + i$ nie powoduje kolizji.
- Najmniej znaczące 0 odpowiada szukanemu, pierwszemu załączkowi, który nie powoduje kolizji.
- Jeśli nie ma 0 w r , sprawdzamy kolejne k załączków...
- ...aż do napotkania bezkolizyjnego albo ostatniego.
- Dla $p(s, c) = (c + \delta s) \bmod L$ i $\delta > 1$, bity odpowiadające nielegalnym przesunięciom wypełniamy jedynkami.



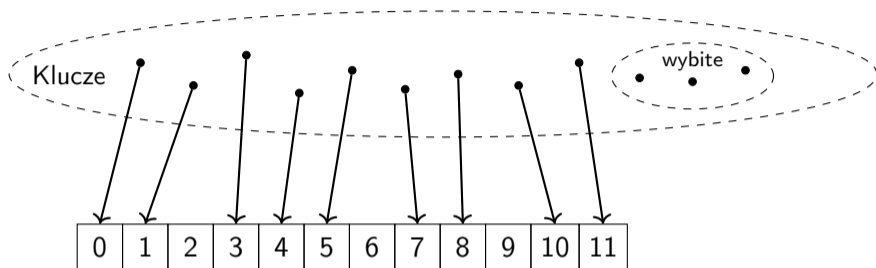
(Minimalne) Doskonałe Funkcje Haszujące



Budowa **doskonałej funkcji haszującej** dla n_1 kluczy polega na skonstruowaniu szeregu funkcji (które mogą wybijać niektóre klucze):

- f_1 mapuje n_1 kluczy wejściowych na $\{0, \dots, n_1 - 1\}$,
- f_2 mapuje wszystkie n_2 klucze wybite przez f_1 na $\{0, \dots, n_2 - 1\}$,
- f_3 mapuje wszystkie n_3 klucze wybite przez f_2 na $\{0, \dots, n_3 - 1\}$,
- ... i tak dalej, aż żadne klucze nie są wybite.

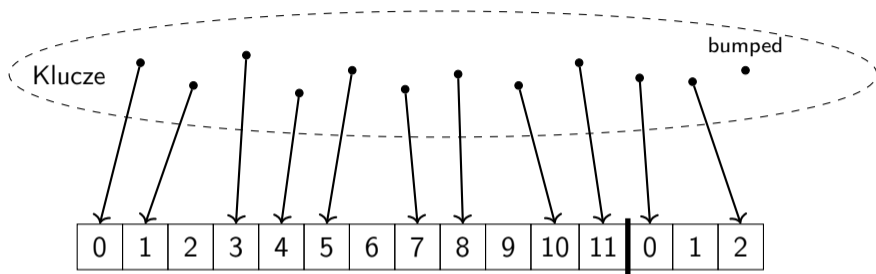
(Minimalne) Doskonałe Funkcje Haszujące



Budowa **doskonałej funkcji haszującej** dla n_1 kluczy polega na skonstruowaniu szeregu funkcji (które mogą wybijać niektóre klucze):

- f_1 mapuje n_1 kluczy wejściowych na $\{0, \dots, n_1 - 1\}$,
- f_2 mapuje wszystkie n_2 klucze wybite przez f_1 na $\{0, \dots, n_2 - 1\}$,
- f_3 mapuje wszystkie n_3 klucze wybite przez f_2 na $\{0, \dots, n_3 - 1\}$,
- ... i tak dalej, aż żadne klucze nie są wybite

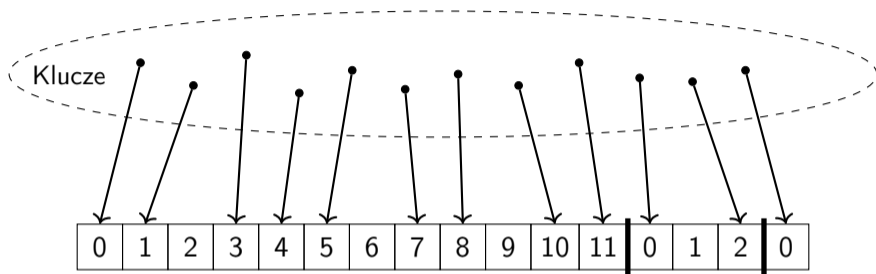
(Minimalne) Doskonałe Funkcje Haszujące



Budowa **doskonałej funkcji haszującej** dla n_1 kluczy polega na skonstruowaniu szeregu funkcji (które mogą wybijać niektóre klucze):

- f_1 mapuje n_1 kluczy wejściowych na $\{0, \dots, n_1 - 1\}$,
- f_2 mapuje wszystkie n_2 klucze wybite przez f_1 na $\{0, \dots, n_2 - 1\}$,
- f_3 mapuje wszystkie n_3 klucze wybite przez f_2 na $\{0, \dots, n_3 - 1\}$,
- ... i tak dalej, aż żadne klucze nie są wybite

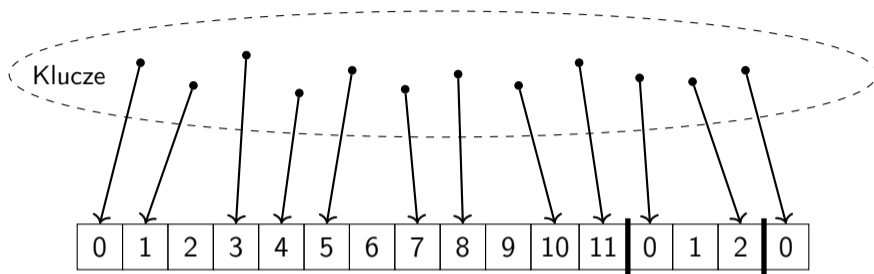
(Minimalne) Doskonałe Funkcje Haszujące



Budowa **doskonałej funkcji haszującej** dla n_1 kluczy polega na skonstruowaniu szeregu funkcji (które mogą wybijać niektóre klucze):

- f_1 mapuje n_1 kluczy wejściowych na $\{0, \dots, n_1 - 1\}$,
- f_2 mapuje wszystkie n_2 klucze wybite przez f_1 na $\{0, \dots, n_2 - 1\}$,
- f_3 mapuje wszystkie n_3 klucze wybite przez f_2 na $\{0, \dots, n_3 - 1\}$,
- ... i tak dalej, aż żadne klucze nie są wybite.

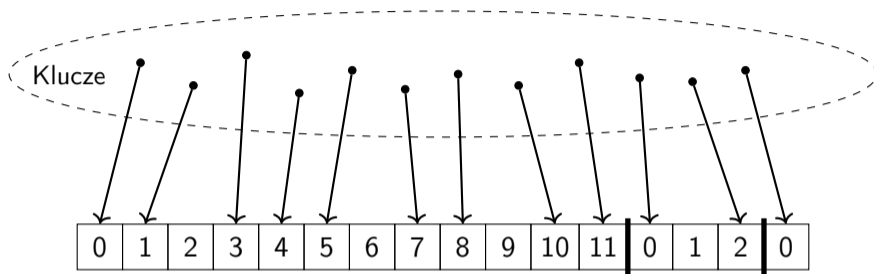
(Minimalne) Doskonałe Funkcje Haszujące



Budowa **doskonałej funkcji haszującej** dla n_1 kluczy polega na skonstruowaniu szeregu funkcji (które mogą wybijać niektóre klucze):

- f_1 mapuje n_1 kluczy wejściowych na $\{0, \dots, n_1 - 1\}$,
- f_2 mapuje wszystkie n_2 klucze wybite przez f_1 na $\{0, \dots, n_2 - 1\}$,
- f_3 mapuje wszystkie n_3 klucze wybite przez f_2 na $\{0, \dots, n_3 - 1\}$,
- ... i tak dalej, aż żadne klucze nie są wybite.

(Minimalne) Doskonałe Funkcje Haszujące

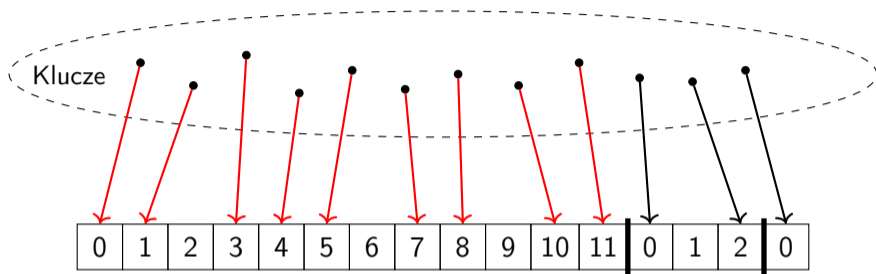


Obliczenie wartości dla klucza k polega na policzeniu $f_1(k), f_2(k), \dots, f_i(k)$, gdzie f_i jest pierwszą funkcją niewybijającą k .

Jeśli już f_1 nie wybija k (co zdarza się dla około $\sim 99\%$ kluczy w typowych konfiguracjach) $f_1(k)$ jest ostatecznie zwróconą wartością.

W przeciwnym razie, jest nią $n_1 + \dots + n_{i-1} + f_i(k)$.

(Minimalne) Doskonałe Funkcje Haszujące

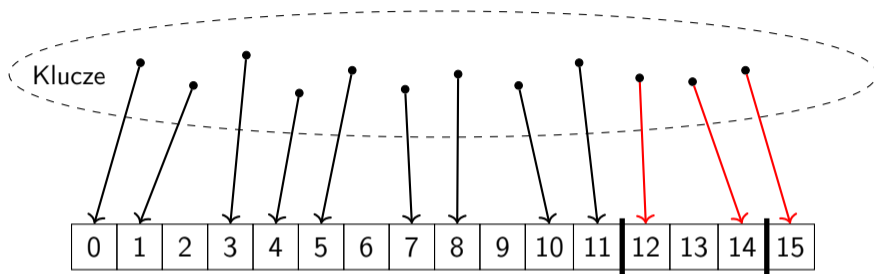


Obliczenie wartości dla klucza k polega na policzeniu $f_1(k), f_2(k), \dots, f_i(k)$, gdzie f_i jest pierwszą funkcją niewybijającą k .

Jeśli już f_1 nie wybija k (co zdarza się dla około $\sim 99\%$ kluczy w typowych konfiguracjach) $f_1(k)$ jest ostatecznie zwróconą wartością.

W przeciwnym razie, jest nią $n_1 + \dots + n_{i-1} + f_i(k)$.

(Minimalne) Doskonałe Funkcje Haszujące

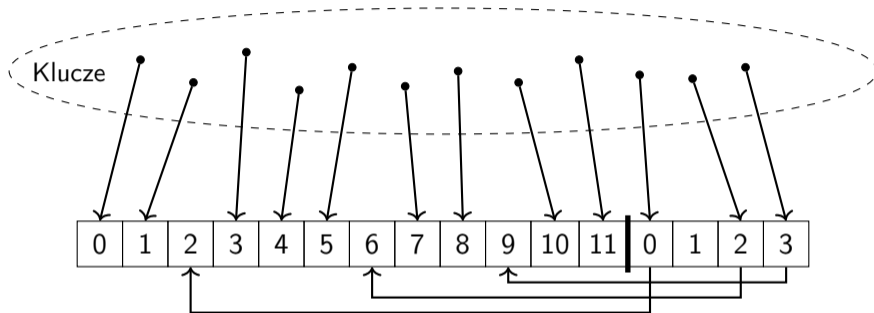


Obliczenie wartości dla klucza k polega na policzeniu $f_1(k), f_2(k), \dots, f_i(k)$, gdzie f_i jest pierwszą funkcją niewybijającą k .

Jeśli już f_1 nie wybija k (co zdarza się dla około $\sim 99\%$ kluczy w typowych konfiguracjach) $f_1(k)$ jest ostatecznie zwróconą wartością.

W przeciwnym razie, jest nią $n_1 + \dots + n_{i-1} + f_i(k)$.

(Minimalne) Doskonałe Funkcje Haszujące

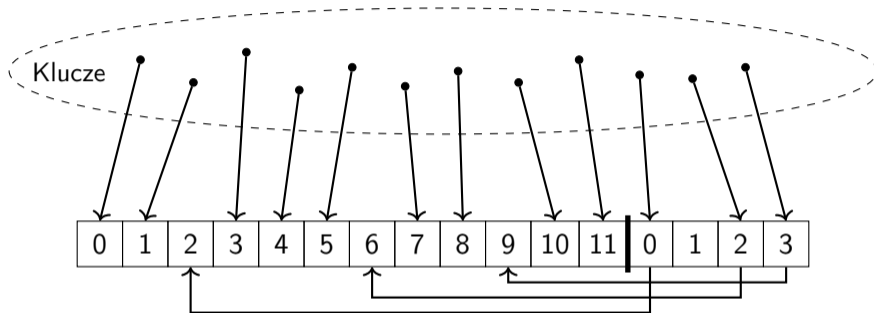


By uzyskać funkcję **minimalną**, używana jest tablica M mapująca kolejne wartości powyżej $n_1 - 1$, na kolejne nieużyte (ze względu na wybijanie przez f_1) wartości poniżej n_1 :

v jest mapowane na $M[v - n_1]$.

M można zwięźle zapisać używając kodowania Elias-Fano.

(Minimalne) Doskonałe Funkcje Haszujące



By uzyskać funkcję **minimalną**, używana jest tablica M mapująca kolejne wartości powyżej $n_1 - 1$, na kolejne nieużyte (ze względu na wybijanie przez f_1) wartości poniżej n_1 :

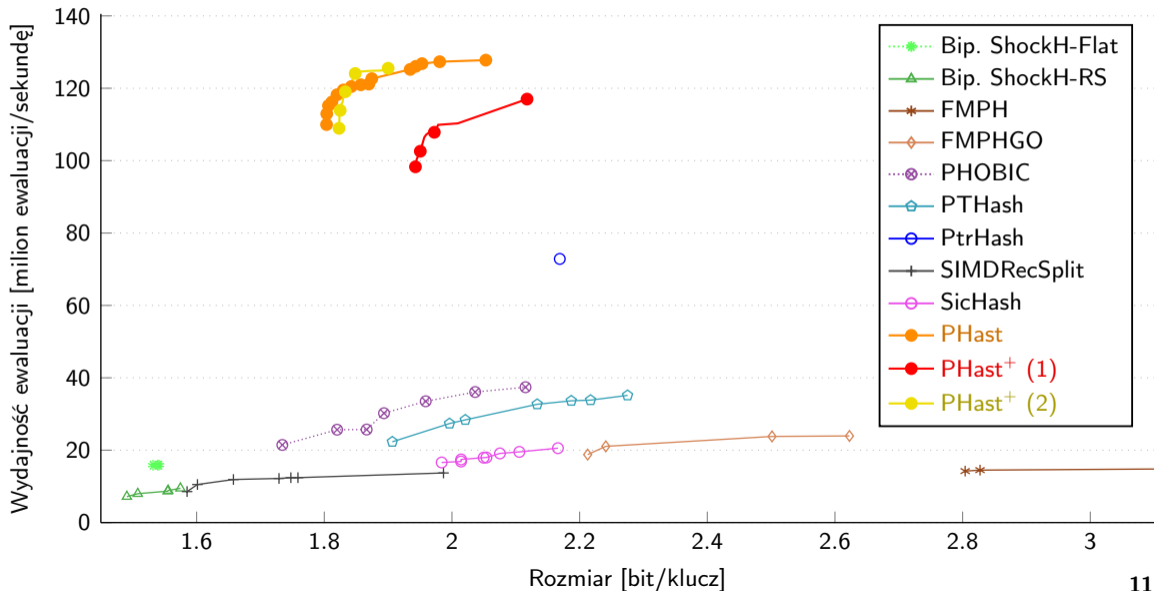
v jest mapowane na $M[v - n_1]$.

M można zwięźle zapisać używając kodowania Elias-Fano.

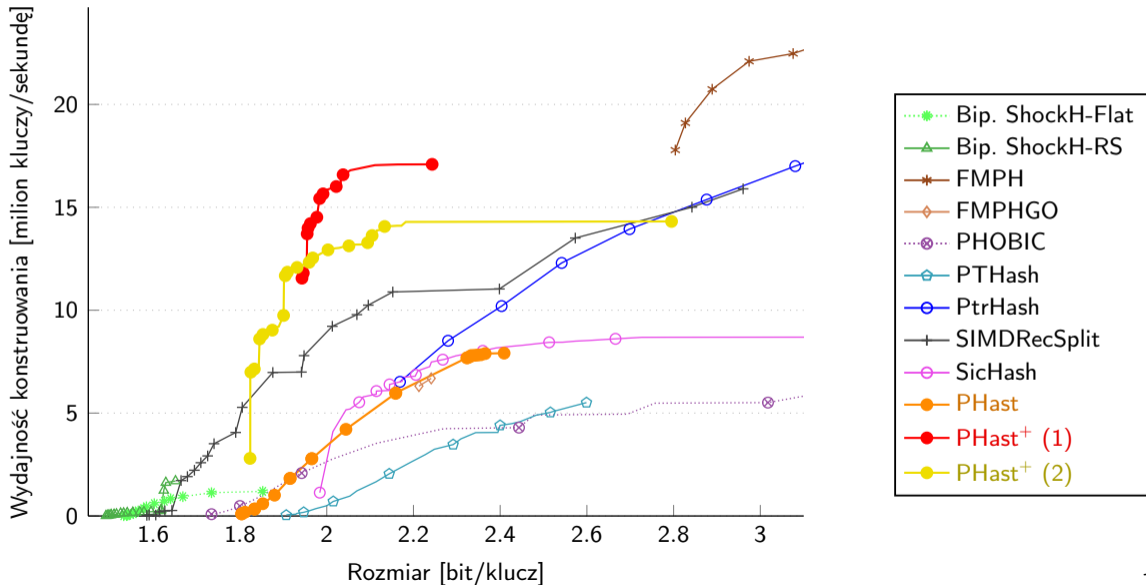
Badanie Wydajności

- badanie za pomocą oprogramowania *MPHF-Experiments* rozwijanego przez Hansa-Petera Lehmana;
- dla 50 milionów kluczy,
- o średniej długości ~ 30 bajtów / klucz;
- konstruowania jedno-wątkowego;
- na AMD Ryzen 5600G @3.9GHz CPU;
- metody zaimplementowane w Rustie (w tym PHast) i C++.

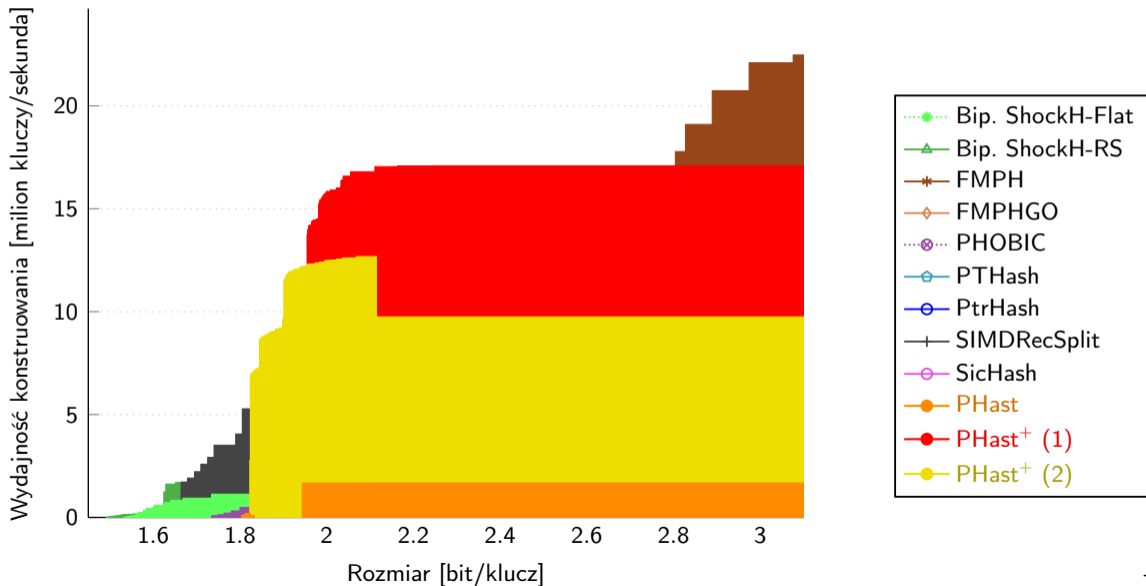
Osiągnięte kompromisy pomiędzy wielkością i szybkością ewaluacji



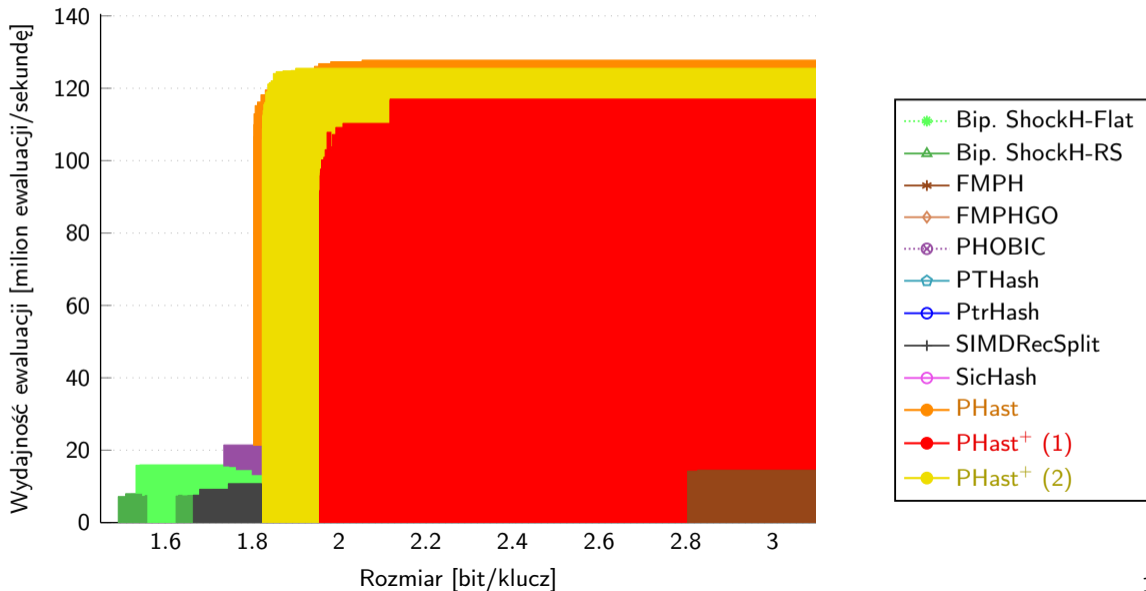
Osiągnięte kompromisy pomiędzy wielkością i szybkością konstruowania



Najszybsze w ewulacji dla zadanych rozmiaru/szybkości konstruowania



Najszybsze do skonstruowania dla zadanych rozmiaru/szybkości ewaluacji



- Piotr Beling, Peter Sanders *PHast - Perfect Hashing made fast* 2026 Proceedings of the SIAM Symposium on Algorithm Engineering and Experiments (ALENEX) pełna, najnowsza wersja dostępna jest na <https://arxiv.org/abs/2504.17918>
- Hans-Peter Lehmann, Thomas Mueller, Rasmus Pagh, Giulio Ermanno Pibiri, Peter Sanders, Sebastiano Vigna, Stefan Walzer, *Modern Minimal Perfect Hashing: A Survey*, ACM Computing Surveys, tom 58, numer 10, dostępny na <https://arxiv.org/abs/2506.06536>
- Stefan Hermann, Hans-Peter Lehmann, Giulio Ermanno Pibiri, Peter Sanders, Stefan Walzer, *PHOBIC: Perfect Hashing With Optimized Bucket Sizes and Interleaved Coding*, ESA 2024
- Giulio Ermanno Pibiri, Roberto Trani, *PTHash: Revisiting FCH Minimal Perfect Hashing*, SIGIR 2021
- Oprogramowanie: <https://crates.io/crates/ph>
<https://github.com/beling/bsuccinct-rs>
<https://github.com/beling/MPHF-Experiments>