

Detecting near-duplicates for web crawling

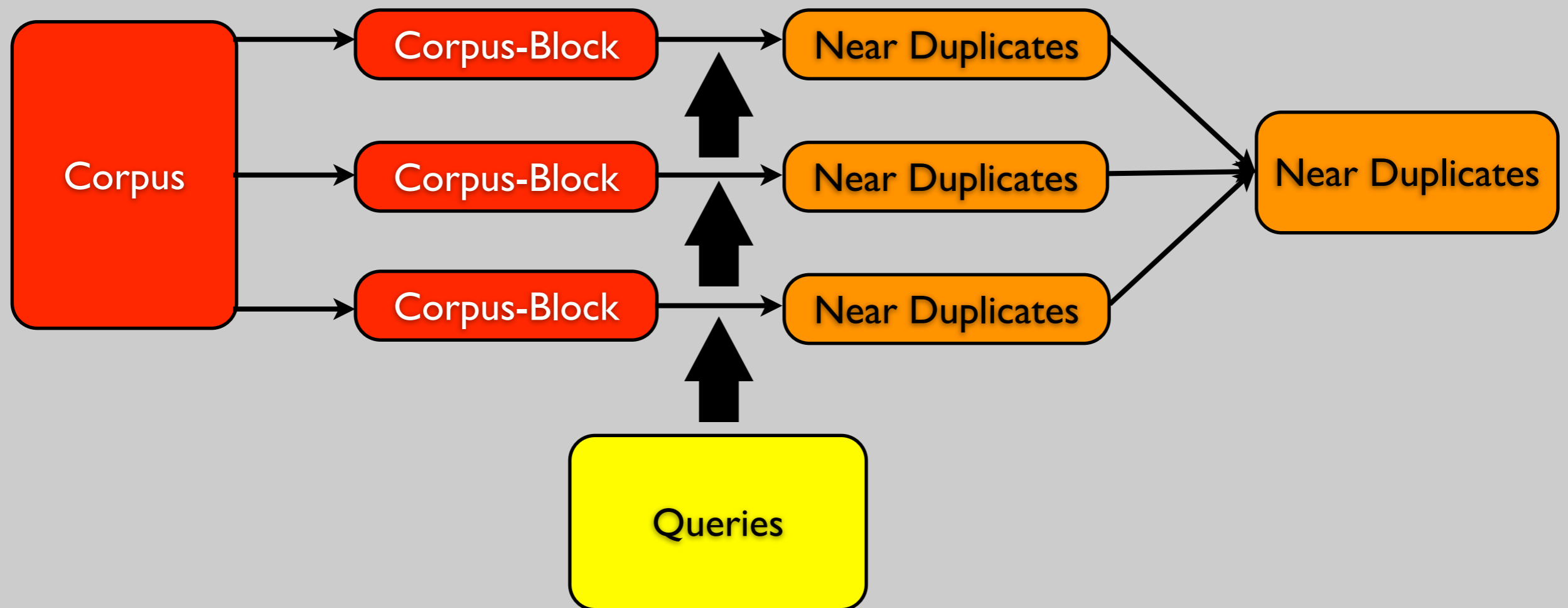
Ziawasch Abedjan
Tobias Flach

Agenda

- Problemstellung
- Generierung der Daten
- Implementierung
- Nächste Schritte

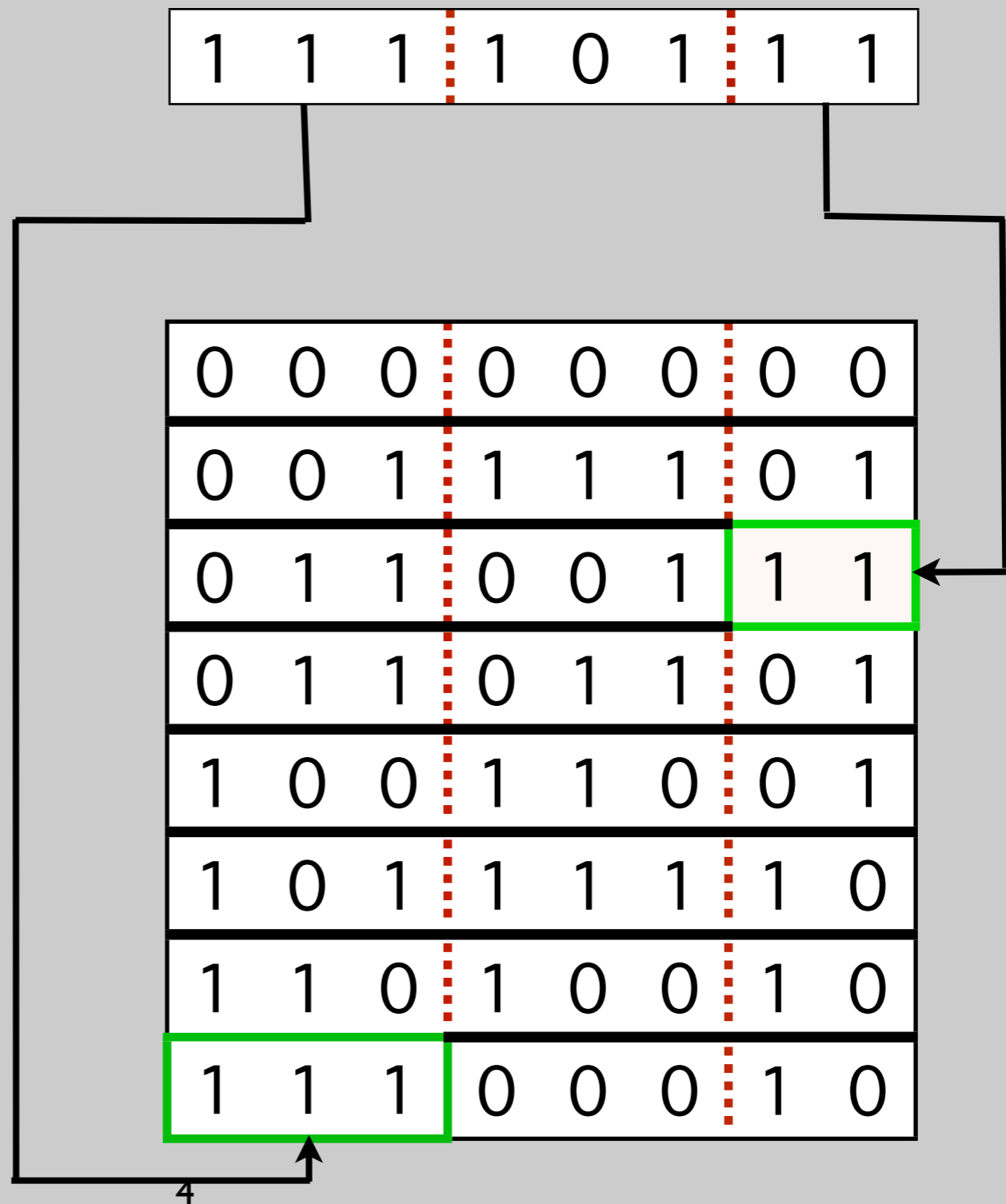
Problemstellung

- Aufgabe: Ähnliche Dokumente finden (near-duplicates)
- Eingabeformat: vorberechnete Hashwerte durch Simhash-Funktion



Vorgehensweise

- Blockaufteilung
- Indizierung
- Identische Suche nach Indexschlüsseln
- Vergleich der Bits



Datengenerierung

- Query-Hashwerte: Zufällige Generierung
- Corpus-Hashwerte:
 - Duplikate: basieren auf Query-Hashwerten, wobei max. k Positionen zufällig invertiert werden
 - Übrige Hashwerte zufällig generiert
- Speicherung als Sequenz von Long-Werten in zwei Dateien

Implementierte Komponenten

- FileInputFormat und RecordReader
 - Einlesen der Daten
 - Aufteilung auf die Mapper
- Map
 - Erzeugung der Masken und Indizes
 - Finden der near-duplicates
- Reduce
 - Verschmelzen der Listen
 - Ausgabe

Finden von near-duplicates

```
List<Long> pdups;
// run over all fingerprints

for(int i = 0; i < fingerprints.length; i++) {
    long fp = fingerprints[i];

    // run over all indices and check for equality for the set mask bits
    for(int j = 0; j < MASKS.length; j++) {
        pdups = indices.get(j).get(MASKS[j] & fp);
        if(pdups != null) {

            // check hamming distance criterion for possible near-duplicates
            for(int k = 0; k < pdups.size(); k++) {
                if(isNearDuplicate(pdups.get(k), fp)) {
                    List<Long> entries = map.get(pdups.get(k));
                    if(!entries.contains(fp)) {
                        entries.add(fp);
                    }
                    map.put(pdups.get(k), entries);
                }
            }
        }
    }
}
```

Finden von near-duplicates

```
public static boolean isNearDuplicate(long x, long y) {
    long hamming_bits = x ^ y;

    int diff = 0;
    for(int i = 0; i < Long.SIZE; i++ ) {
        diff += hamming_bits & 1;
        hamming_bits >>= 1;
        if(diff > MAX_DISTANCE) {
            return false;
        }
    }

    return true;
}
```


Ausgabe

File: [/user/toby/out/part-00000](#)

Goto :

[Go back to dir listing](#)

[Advanced view/download options](#)

-7253127574651717233	-7253126475140092529	-7253127849529755250	-2641441556223280757
-5491288702235061828	-5491288702235324004	-6639706607214554692	-5491288702503496708
-5486317444252353411	-4909293741995508483	-5487144276996440963	-5486317444244161411
-4797612092326257425	-4815626490835739161	-4793099696605864705	-5085842468478223121
-4180025748777018677	-4180025748743456565	-4180025748772757813	-4184529485977560373
-3544617457957310509	-3544617440240570477	-3544615258866946095	-3544617457424633902
-3222778774478489740	-3222778774478490288	-3222778774461646988	-3223904622845724812
-1894294158393768186	-1894294158670723322	-1750178970319026426	-1930323024131684602
-1847393393625710773	-1802357397352006325	-1847393410806104241	-1847428577997801621
63805910251109489	63805910183998585	189906562378529905	1216727346138479729
409528470880629098	409528473027064171	2715362684001300586	-8813843565907562134
3070591683186794097	1917670178311511665	3070590858553073249	764748673975196785
4325293537216136632	4327543138006566296	4325293520170493368	3892947972988569528
5042447885026469617	5042446802761819889	5042447340639364849	5042447884892316401
5752902165093557664	5753042902580868512	5734887732224341408	5752902714580936097
7724805193449252538	7724805193449218738	7724805193583462074	7760833991005152954
7877711739454245429	7916133073775249973	7877746958194461237	7877746924094769717
8257354993143999154	8545580971248937650	8275369391653481140	8221607668954262194
8388387801348499726	8370373471558495502	8388387526474788110	8964848553651923338
8559203500653729183	9135664252957119967	8559203432001099167	8505177897311327647

Nächste Schritte

- Korrektur des RecordReaders
- Performance-Analyse des Algorithmus
- Testen der Abhängigkeit von der Werteverteilung (u.a. ähnliche Hashwerte, Einfluss einer hohen Hamming Distance)
- Mögliche Erweiterungen:
 - Verteilung der Indizes auf verschiedene Nodes (mit anschließender Duplikatentfernung)
 - Verteilte Berechnung des Simhash-Wertes