

Addendum to *Scalable Splitting of Massive Data Streams*

Erik Zeitler, Tore Risch

Department of Information Technology
Uppsala University
Sweden
erik.zeitler@it.uu.se
tore.risch@it.uu.se

Abstract. Equation (9) in Zeitler, Risch: *Scalable Splitting of Massive Data Streams* (Proc. DASFAA 2010) described how to compute the fanout of each level in a theoretically optimal splitstream tree, called *maxtree*. This addendum shows how the *maxtree* formula is derived.

1 Derivation of the *maxtree* formula

The cost model for a node at level ℓ in a splitstream tree is expressed using Equation (5) in [1]:

$$C_\ell = \Phi o^{(\ell-1)} \cdot (cc + (cp + ce) \cdot (r_\ell + f_\ell \cdot b_\ell)) \quad (1)$$

Using this model, the cost at the root node C_1 (whose fanout is $f_1 = 2$) is

$$C_1 = \Phi \cdot (cc + (cp + ce) \cdot (r + 2 \cdot b)) \quad (2)$$

As discussed in [1], solving f_ℓ for $C_\ell = C_1$ gives the maximum allowed (optimal) fanout at each level $\ell > 1$.

$$\begin{aligned} \Phi o^{(\ell-1)} \cdot (cc + (cp + ce) \cdot (r_\ell + f_\ell \cdot b_\ell)) &= \\ \Phi \cdot (cc + (cp + ce) \cdot (r + 2 \cdot b)) & \end{aligned} \quad (3)$$

Substitute r_ℓ , b_ℓ , and $\Phi o^{(\ell)}$ using Equations (4) and (5) in [1] to express Equation (3) in terms of r , b , and $\lambda_{\ell-1}$:

$$\begin{aligned} \Phi \left(b + \frac{r}{\lambda_{\ell-1}} \right) \cdot \left(cc + (cp + ce) \cdot \left(\frac{r + f_\ell \cdot b \lambda_{\ell-1}}{r + b \lambda_{\ell-1}} \right) \right) &= \\ \Phi \cdot (cc + (cp + ce) \cdot (r + 2 \cdot b)) & \end{aligned} \quad (4)$$

It is easy to see that Φ cancels. Multiply both sides by $\lambda_{\ell-1} / (cp + ce)$:

$$(r + \lambda_{\ell-1}b) \cdot \left(\frac{cc}{cp + ce} + \left(\frac{r + f_\ell \cdot \lambda_{\ell-1}b}{r + \lambda_{\ell-1}b} \right) \right) = \quad (5)$$

$$\lambda_{\ell-1} \frac{cc}{cp + ce} + \lambda_{\ell-1}r + 2 \cdot \lambda_{\ell-1}b$$

Set $a = cc / (cp+ce)$ and simplify:

$$f_\ell \cdot \lambda_{\ell-1}b = \quad (6)$$

$$a(\lambda_{\ell-1} - r - \lambda_{\ell-1}b) + \lambda_{\ell-1}r + 2 \cdot \lambda_{\ell-1}b - r$$

Divide by $\lambda_{\ell-1}b$:

$$f_\ell = \frac{a}{b} - \frac{ar}{\lambda_{\ell-1}b} - a + \frac{r}{b} + 2 - \frac{r}{\lambda_{\ell-1}b} \quad (7)$$

$$= 2 + \frac{r}{b} \left(1 - \frac{a}{\lambda_{\ell-1}} - \frac{1}{\lambda_{\ell-1}} \right) + \frac{a}{b}(1-b)$$

As each tuple is routed, broadcasted, or omitted, $r + b + o = 1$. No tuples are omitted in the splitstream tree ($o = 0$), as assumed in the beginning of Section 3.2 of [1]. Thus, $r = 1 - b$, which is used to obtain the final formula:

$$f_\ell = 2 + \frac{r}{b} \left(1 - \frac{a}{\lambda_{\ell-1}} - \frac{1}{\lambda_{\ell-1}} \right) + \frac{ar}{b} \quad (8)$$

$$= 2 + \frac{r}{b} \left(1 + a - \frac{a}{\lambda_{\ell-1}} - \frac{1}{\lambda_{\ell-1}} \right)$$

$$= 2 + \frac{r}{b} (1 + a) \left(1 - \frac{1}{\lambda_{\ell-1}} \right)$$

maxtree is constructed by first assigning two child nodes to the root node. Then, for each node at level ℓ , add f_ℓ child nodes. Keep adding levels until λ_ℓ is at least w .

2 References

1. Zeitler, E., Risch, T.: Scalable Splitting of Massive Data Streams. In: DASFAA (2010)