

# $\alpha$ - $\beta$ -Factorisation and the Binary Case of Simon's Congruence

Pamela Fleischmann<sup>(A)</sup>   Jonas Höfer<sup>(B)</sup>   Annika Huch<sup>(A)</sup>  
Dirk Nowotka<sup>(A)</sup>

<sup>(A)</sup>Kiel University, Kiel, Germany

fpa@informatik.uni-kiel.de, stu216885@mail.uni-kiel.de, dn@informatik.uni-kiel.de

<sup>(B)</sup>University of Gothenburg, Sweden

jonas.hofer@gu.se

## Abstract

Based on the arch factorisation (Hébrard 1991), first the notion of  $k$ -richness and later the one of  $k$ -universality - both measure words by their scattered factors - were introduced. In 2022 Fleischmann et al. presented a generalisation by intersecting the arch factorisations of a word and its reverse. Here we use this  $\alpha$ - $\beta$ -factorisation in order to characterise the Simon congruence of  $k$ -universal words in terms of 1-universal words and apply these results to binary words obtaining a full characterisation of the index of the congruence.

## 1. Introduction

A *scattered factor*, *subsequence*, or *scattered subword* of a word  $w$  is a word that is obtained by deleting letters from  $w$  while preserving the order of the remaining ones, e.g., *tea* and *thora* are both scattered factors of *theorietag*. In contrast to a factor, like *eta*, a scattered factor is not necessarily contiguous. Here, we focus on Simon's congruence [8]  $\sim_k$  for  $k \in \mathbb{N}_0$ :  $u \sim_k v$  iff  $u, v$  share all scattered factors up to length  $k$ . A long outstanding question, posed by Sakarovitch and Simon [7], is the exact structure of the congruence classes of  $\sim_k$  and the index of the relation. Currently, no exact formula is known. One approach for studying scattered factors in words is based on the notion of *scattered factor universality* [1, 2, 3]. A word  $w$  is called  $\ell$ -*universal* if it contains all words of length  $\ell$  as scattered factors. For instance, the word *alfalfa*<sup>1</sup> is 2-universal since it contains all words of length two over the alphabet  $\{a, l, f\}$  as scattered factors. A main tool in this line of research is the  $\alpha$ - $\beta$ -*factorization* [3]. Kosche et al. [6] implicitly used this factorisation to determine shortest absent scattered factors in words.

*Our Contribution.* We investigate the  $\alpha$ - $\beta$ -factorization and give necessary and sufficient conditions for the congruence of words in terms of their factors. We characterise  $\sim_k$  in terms of 1-universal words through their  $\alpha\beta\alpha$ -factors. We use these results to characterize the classes of binary words and their cardinality, as well as, to calculate the index in this special case. Lastly, we start to transfer the previous results to the ternary alphabet.

---

<sup>1</sup>Alfalfa (*Medicago sativa*) is plant whose name means *horse food* in Old Persian

## 2. Preliminaries

Let  $\mathbb{N} = \{1, 2, \dots\}$  and set  $\mathbb{N}_0 = \{0\} \cup \mathbb{N}$ ,  $[m] = \{1, \dots, m\}$ , and  $[m]_0 = \{0\} \cup [m]$ . For the standard definitions of combinatorics on words, we refer to [7]. We abbreviate an alphabet of cardinality  $i \in \mathbb{N}$  by  $\Sigma_i$ . If  $w = xy$  we write  $x^{-1}w$  for  $y$  and  $wy^{-1}$  for  $x$ . A word  $u \in \Sigma^*$  of length  $n \in \mathbb{N}_0$  is called a *scattered factor* of  $w \in \Sigma^*$  if there exist  $v_0, \dots, v_n \in \Sigma^*$  with  $w = v_0u[1]v_1 \cdots v_{n-1}u[n]v_n$ . Let  $\text{ScatFact}(w)$ ,  $\text{ScatFact}_k$ , and  $\text{ScatFact}_{\leq k}$  denote the sets of all, exactly of length  $k$ , up to length  $k$  resp. scattered factors of  $w$ . For comparing words w.r.t. their scattered factors, Simon introduced a congruence relation nowadays known as *Simon's congruence* [8]: two words  $u, v \in \Sigma^*$  are called *Simon  $k$ -congruent* ( $u \sim_k v$ ) iff  $\text{ScatFact}_{\leq k}(u) = \text{ScatFact}_{\leq k}(v)$  for some  $k \in \mathbb{N}$ . A word  $w \in \Sigma^*$  is called  *$k$ -universal* w.r.t.  $\Sigma$  if  $\text{ScatFact}_k(w) = \Sigma^k$ . The maximal  $k$  such that  $w$  is  $k$ -universal is denoted by  $\iota(w)$  and called  $w$ 's *universality index*. For a word  $w \in \Sigma^*$  the *arch factorisation* is given by  $w = \text{ar}_1(w) \cdots \text{ar}_k(w) \text{re}(w)$  for  $k \in \mathbb{N}_0$  with  $\text{alph}(\text{ar}_i(w)) = \Sigma$  for all  $i \in [k]$ , the last letter of  $\text{ar}_i(w)$  occurs exactly once in  $\text{ar}_i(w)$  for all  $i \in [k]$ , and  $\text{alph}(\text{re}(w)) \subset \Sigma$ . The words  $\text{ar}_i(w)$  are called *arches* and  $\text{re}(w)$  is the *rest* of  $w$ . Define the *modus* of  $w$  as  $\text{m}(w) = \text{ar}_1(w)[|\text{ar}_1(w)|] \cdots \text{ar}_k(w)[|\text{ar}_k(w)|] \in \Sigma^k$ . Set  $\text{ar}_{i..j}(w) = \text{ar}_i(w) \cdots \text{ar}_j(w)$ . The  $\alpha$ - $\beta$ -factorisation was introduced in [3] inspired by [6]. Define for the arch factorisation of  $w^R$  (read left to right) the  $i^{\text{th}}$  *reverse arch*  $\tilde{\text{ar}}_i(w) = (\text{ar}_{\iota(w)-i+1}(w^R))^R$ , the *reverse rest*  $\tilde{\text{re}}(w) = (\text{re}(w^R))^R$ , and set  $\tilde{\text{m}}(w)$  as  $\text{m}(w^R)^R$  for the *reverse modus*.

**Definition 2.1** For  $w \in \Sigma^*$  define  $w$ 's  $\alpha$ - $\beta$ -factorisation by  $w = \alpha_0\beta_1\alpha_1 \cdots \alpha_{\iota(w)-1}\beta_{\iota(w)}\alpha_{\iota(w)}$  with  $\text{ar}_i(w) = \alpha_{i-1}\beta_i$  and  $\tilde{\text{ar}}_i(w) = \beta_i\alpha_i$  for all  $i \in [\iota(w)]$ ,  $\tilde{\text{re}}(w) = \alpha_0$ , as well as  $\text{re}(w) = \alpha_{\iota(w)}$ . Define  $\text{core}_i = \beta_i[2..|\beta_i| - 1]$  or  $\varepsilon$  if  $|\beta_i| \leq 2$ .

Note that the  $\alpha$ - $\beta$ -factorisation is left-right-symmetric and that the  $i^{\text{th}}$  reverse arch always starts inside the  $i^{\text{th}}$  arch. We finish this section with three results from [4, 5, 8]

**Lemma 2.2** Let  $u, v \in \Sigma^*$ ,  $u', v' \in \Sigma^+$ , and  $\mathbf{x} \in \Sigma$ .

- (1) If  $u \sim_k v$  then  $w_1uw_2 \sim_{\iota(w_1)+k+\iota(w_2)} w_1vw_2$ .
- (2)  $u'v' \sim_k u'$  iff  $u' = u'_1 \cdots u'_k$  such that  $\text{alph}(u'_1) \supseteq \dots \supseteq \text{alph}(u'_k) \supseteq \text{alph}(v')$ .
- (3)  $wv \sim_k uxv$  iff there exist  $p, p' \in \mathbb{N}_0$  with  $p + p' \geq k$  and  $ux \sim_p u$  and  $xv \sim_{p'} v$ .

## 3. $\alpha$ - $\beta$ -Factorisation

In this section, we investigate the  $\alpha$ - $\beta$ -factorisation based on results of [4]. The main result states that it suffices to look at 1-universal words in order to gain the information about the  $\sim_k$  congruence classes. First, we show that *cutting of  $\ell$  arches* from two  $k$ -congruent words each, leads to  $(k - \ell)$ -congruence. Then we connect the congruence of words to the congruence of their  $\alpha$  factors, leading to a characterisation by the congruence of  $\alpha\beta\alpha$ -factors.

**Lemma 3.1** Let  $w, \tilde{w} \in \Sigma^*$  with  $w \sim_k \tilde{w}$  and  $\iota(w) = \iota(\tilde{w}) < k$ , then  $\text{ar}_1^{-1}(w) \cdot w \sim_{k-1} \text{ar}_1^{-1}(\tilde{w}) \cdot \tilde{w}$  and  $\alpha_i\beta_{i+1}\alpha_{i+1} \cdots \alpha_j \sim_{k-\iota(w)+j-i} \tilde{\alpha}_i\tilde{\beta}_{i+1}\tilde{\alpha}_{i+1} \cdots \tilde{\alpha}_j$  for all  $0 \leq i \leq j \leq \iota(w)$ .

**Proposition 3.2** For all  $w, \tilde{w} \in \Sigma^*$  with  $m = \iota(w) = \iota(\tilde{w}) < k$  such that  $\beta_i = \tilde{\beta}_i$  for all  $i \in [m]$ , we have  $w \sim_k \tilde{w}$  iff  $\alpha_i \sim_{k-m} \tilde{\alpha}_i$  for all  $i \in [m]_0$ . Thus,  $w \sim_k \tilde{w}$  iff  $\alpha_i \sim_{k-m} \tilde{\alpha}_i$  for all  $i \in [m]_0$  and for  $w' = \alpha_0\tilde{\beta}_1\alpha_1 \cdots \tilde{\beta}_m\alpha_m$  we have  $w \sim_k w'$ .

**Theorem 3.3** *Let  $w, \tilde{w} \in \Sigma^*$  with  $m = \iota(w) = \iota(\tilde{w}) < k$ . Then,  $w \sim_k \tilde{w}$  iff  $\alpha_{i-1}\beta_i\alpha_i \sim_{k-m+1} \tilde{\alpha}_{i-1}\tilde{\beta}_i\tilde{\alpha}_i$  for all  $i \in [m]$ .*

In the light of Theorem 3.3, in the following, we consider some special cases of these triples w.r.t. the alphabet of the both involved  $\alpha$ . Hence, let  $w, \tilde{w} \in \Sigma^*$  with  $1 = \iota(w) = \iota(\tilde{w})$ .

**Proposition 3.4** (1) *Let  $\alpha_0 = \alpha_1 = \tilde{\alpha}_0 = \tilde{\alpha}_1 = \varepsilon$ . Then  $w \sim_k \tilde{w}$  iff  $k = 1$  or  $k \geq 2$ ,  $m(w) = m(\tilde{w})$ ,  $\tilde{m}(w) = \tilde{m}(\tilde{w})$ , and  $\text{core}_1 \sim_k \widetilde{\text{core}}_1$ .*

(2) *Let  $\text{alph}(\alpha_i) = \text{alph}(\tilde{\alpha}_i) \in \binom{\Sigma}{|\Sigma|-1}$ , then  $w \sim_k \tilde{w}$  iff  $\alpha_i \sim_{k-1} \tilde{\alpha}_i$  for all  $i \in [1]_0$ .*

The last proposition does not hold if not both  $m(w)$  and  $\tilde{m}(w)$  are identical: consider  $w = \text{ababeabab} \cdot \text{abcd} \cdot \text{cdcdcd} \sim_4 \text{ababeabab} \cdot \text{baedc} \cdot \text{cdcdcd} = \tilde{w}$  with  $m(w) = d \neq c = m(\tilde{w})$  and  $\tilde{m}(w) = a \neq b = \tilde{m}(\tilde{w})$ . In the next proposition, we give a necessary condition for the  $\alpha$ -factors.

**Proposition 3.5** *Let  $w \in \Sigma^*$  with  $\iota(w) = 1$ ,  $k \in \mathbb{N}$ , and  $\tilde{M} = \{\tilde{m}(\tilde{w})[1] \mid \tilde{w} \in [w]_{\sim_k}\}$ . If  $|\tilde{M}| \geq 2$  then there exists a factorisation  $\alpha_0 = u_1 \cdots u_{k-1}$  with  $\text{alph}(u_1) \supseteq \dots \supseteq \text{alph}(u_{k-1}) \supseteq \tilde{M}$ .*

## 4. The Binary and Ternary Case of Simon's Congruence

First, we apply our results to the binary alphabet. Note that for a given  $w$  with  $\iota(w) \leq k$ , we have  $|\{\tilde{m}(\tilde{w}) \mid \tilde{w} \in [w]_{\sim_k}\}| = 1$ .

**Proposition 4.1** *For all  $w \in \Sigma_2^*$ , we have for all  $i \in [\iota(w)]$ ,  $\beta_i \in \{a, b, ab, ba\}$ . If  $\beta_i = x$ , then  $\alpha_{i-1}, \alpha_i \in \bar{x}^+$  with  $x \in \Sigma_2$  and if  $\beta_i = x\bar{x}$ , then  $\alpha_{i-1} \in x^*$  and  $\alpha_i \in \bar{x}^*$  with  $x \in \Sigma_2$ .*

In the binary case the  $k$ -congruence of two words with identical  $\iota < k$  leads to the same modi and same  $\beta$  giving a characterisation of  $\sim_k$  for binary words in terms of unary words.

**Lemma 4.2** *Let  $w, w' \in \Sigma_2^*$  with  $w \sim_k w'$  and  $m = \iota(w) = \iota(w') < k$ , then  $m(w) = m(w')$  and thus,  $\beta_i = \beta'_i$  for all  $i \in [m]$ .*

**Theorem 4.3** *Let  $w, w' \in \Sigma_2^*$  such that  $m = \iota(w) = \iota(w') < k$ , then  $w \sim_k w'$  iff  $\beta_i = \beta'_i$  for all  $i \in [m]$  and  $\alpha_i \sim_{k-m} \alpha'_i$  for all  $i \in [m]_0$ .*

Theorem 4.3 implies if  $|[w]_{\sim_k}| = \infty$ , then  $x^k \in \text{ScatFact}_k(w)$  for some  $x \in \Sigma$  (the contrary is generally not true:  $v = \text{bbabb}$  w.r.t.  $\sim_4$ ). The following theorem characterises the binary case.

**Theorem 4.4** *Let  $w \in \Sigma_2^*$ , then  $|[w]_{\sim_k}| < \infty$ . We have  $|[w]_{\sim_k}| = 1$  iff  $\iota(w) < k$  and  $|\alpha_i| < k - \iota(w)$  for all  $i \in [\iota(w)]_0$ .*

We present a formula for the precise value of  $|\Sigma_2^*/\sim_k|$  counting classes based on the valid combinations of  $\beta$ -factors and number of classes for each  $\alpha$ -factors and giving the index.

**Theorem 4.5** *The number of congruence classes of  $\Sigma_2^*/\sim_k$  of words with  $m < k$  arches is given by  $\left\| \binom{k-m}{k-m} \binom{k-m}{k-m} \binom{k-m}{k-m} \cdot \binom{k-m}{k-m} \right\|_1 = c_k^m$  where  $c_k^{-1} = 1$ ,  $c_k^0 = 2k + 1$ , and  $c_k^m = 2 \cdot (k - m + 1) \cdot c_{k-1}^{m-1} - 2 \cdot (k - m) \cdot c_{k-2}^{m-2}$  where  $\|\cdot\|_1$  denotes the 1-norm.*

**Corollary 4.6** For  $k \in \mathbb{N}_0$  we have  $|\Sigma_2^*/\sim_k| = 1 + \sum_{m=0}^{k-1} c_k^m$ .

Now, we consider  $\Sigma_3$ . Note that if  $m_1(w) = \tilde{m}_1(w)$  then  $\text{core}_1 = \varepsilon$ . If  $m_1(w) \neq \tilde{m}_1(w)$ , then  $\text{core}_1 \in (\Sigma \setminus \{m_1(w), \tilde{m}_1(w)\})^*$ , i.e., the cores are unary and denoted by  $y \in \Sigma_3$ . Define for a boolean predicate  $P$ ,  $\delta_{P(x)} = 1$  if  $P(x)$  is true and 0 otherwise. We assume  $k \geq 2$  (we characterise 1-universal words) and  $w, \tilde{w} \in \Sigma_3^*$  with  $1 = \iota(w) = \iota(\tilde{w})$ .

**Lemma 4.7** Let  $m(w) = m(\tilde{w})$  and  $\tilde{m}(w) = \tilde{m}(\tilde{w})$ , we have  $w \sim_k \tilde{w}$  iff  $\alpha_i \sim_{k-1} \tilde{\alpha}_i$  for all  $i \in [1]_0$  and  $\text{core}_1 \sim_{k-c} \widetilde{\text{core}}_1 \in y^*$  where  $c := \iota(\alpha_0) + \delta_{y \preceq \text{re}(\alpha_0)} + \iota(\alpha_1) + \delta_{y \preceq \tilde{\text{re}}(\alpha_1)}$ .

We finish with a characterisation in the ternary case.

**Theorem 4.8** For  $w, \tilde{w} \in \Sigma_3^*$  we have  $w \sim_k \tilde{w}$  iff  $\alpha_i \sim_{k-1} \tilde{\alpha}_i$  for all  $i \in [1]_0$ , and  
(1)  $|\text{alph}(\alpha_i)| = 2$ ,  $\text{alph}(\alpha_{1-i}) \cap \text{alph}(\alpha_i) = \emptyset$ , and  $\iota(\alpha_i) \geq k-1$  for some  $i \in [1]_0$ , or  
(2)  $m(w) = m(\tilde{w})$ ,  $\tilde{m}(w) = \tilde{m}(\tilde{w})$ , and  $\text{core} \sim_{k-c} \widetilde{\text{core}}$  where  $c := \iota(\alpha_0) + \delta_{y \preceq \alpha_0} + \iota(\alpha_1) + \delta_{y \in \alpha_1}$ .

## 5. Conclusion

In this paper, we investigated the  $\alpha$ - $\beta$ -factorisation (cf. [6, 3]) as an object of intrinsic interest. This leads to a result characterising  $k$ -congruence of  $m$ -universal words in terms of their 1-universal  $\alpha\beta\alpha$ -factors. In the case of the binary and ternary alphabet, we fully characterised the congruence of words in terms of their single factors. Extending this idea of the  $\alpha$ - $\beta$ -factorisation to lower layers (arches w.r.t. some  $\Omega \subset \Sigma$ ), is left as future work.

## References

- [1] L. BARKER, P. FLEISCHMANN, K. HARWARDT, F. MANEA, D. NOWOTKA, Scattered factor-universality of words. In: *DLT*. Springer, 2020, 14–28.
- [2] P. FLEISCHMANN, S. GERMANN, D. NOWOTKA, Scattered Factor Universality–The Power of the Remainder. *preprint arXiv:2104.09063 (published at RuFiDim) (2021)*.
- [3] P. FLEISCHMANN, L. HASCHKE, A. HUCH, A. MAYROCK, D. NOWOTKA, Nearly  $k$ -universal words–investigating a part of simon’s congruence. In: *DCFS*. 2022, 57–71.
- [4] P. KARANDIKAR, M. KUFLEITNER, P. SCHNOEBELEN, On the index of Simon’s congruence for piecewise testability. *Inf. Process. Lett.* **115** (2015) 4, 515–519.
- [5] P. KARANDIKAR, P. SCHNOEBELEN, The height of piecewise-testable languages and the complexity of the logic of subwords. *LICS* **15** (2019) 2.
- [6] M. KOSCHE, T. KOSS, F. MANEA, S. SIEMER, Absent subsequences in words. In: *RP*. Springer, 2021, 115–131.
- [7] M. LOTHAIRE, *Combinatorics on Words*. Cambridge Mathematical Library, Cambridge University Press, 1997.
- [8] I. SIMON, Piecewise testable events. In: *Autom. Theor. Form. Lang., 2nd GI Conf.*. LNCS 33, Springer, 1975, 214–222.