

# Traitements des similarités dans la justification de paragraphe

Exposé GUTenberg, 10 septembre 2025

Didier Verna

[didier@didierverna.net](mailto:didier@didierverna.net)



[didierverna.net](http://didierverna.net)



@didierverna



didier.verna



in/didierverna

# Préambule

## ► Article à ACM SIGWEB DocEng 2024

### Similarity Problems in Paragraph Justification

An Extension to the Knuth-Plass Algorithm

Didier Verna  
EPITA Research Laboratory  
Le Kremlin-Bicêtre, France  
didier@lefr.epita.fr

#### Abstract

In high quality typography, consecutive lines beginning or ending with the same word or sequence of characters is considered a defect. We have implemented an extension to TeX's paragraph justification algorithm which handles this problem. Experimentation shows that given a set of rules, our extension can detect and avoid such cases. Our extension automates the detection and avoidance of similarities while leaving the ultimate decision to the professional typographer, thanks to a new adjustable curve. The extension is simple and lightweight, making it a useful addition to production engines.

#### CCS Concepts

- Applied computing → Document preparation; - Theory of computation → Dynamic graph algorithms.

#### Keywords

Paragraph Justification, Similarity Avoidance, Homoeoteleutons, Homoiotropes, RGK, Knuth-Plass Extension

#### ACM Reference Format

Verna, D. 2024. Similarity Problems in Paragraph Justification: An Extension to the Knuth-Plass Algorithm. In *ACM Symposium on Document Engineering 2024 (DocEng '24)*. August 20–25, 2024, San Jose, CA, USA. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3685408.3685666>

#### 1 Introduction

In spite of its relatively old age, Donald Knuth's TeX typesetting system [10, 11] is still considered a *de facto* standard when it comes to digital typography. In particular, its paragraph justification algorithm, known as the Acute Class [12], established a landmark in the catalog of algorithms combining a paragraph as a whole rather than proceeding line by line, as earlier (greedy) algorithms used to do [1, 4, 16].

Yet, many aspects of fine typography are not directly or automatically handled by TeX. Consider for example the leftmost paragraph in Figure 1. It is justified by line, and by line, in a 10pt font, a Roman font at a 10pt size, and for a paragraph width of 201pt (the figure is optically scaled down in order to fit on the page).

Notice how these lines near the end of the paragraph and the same way, with the word "and". This is considered a defect in high quality typography, as it generates a micro-interruption: the reader's attention may be caught by the similarity and temporarily diverted from the main text flow. Such a defect may also lead the reader to accidentally skip a letter to read the next one, even more so when the problem occurs at the beginning of the line rather than at the end of it. In fact, the field of textual criticism has identified the problem and its consequences in the very ancient context of scroll errors (e.g. missing lines in manual copies of the bible made by monks) [20].

We have implemented an extension to the Knuth-Plass (KP) algorithm that is able to deal with that kind of defect. In this paper, we use the term "similarity" the lack of nice French typography. We would like to prevent the same character "word key" analogous to the more widespread expression *hyphenation ladder*, the meaning of which should be obvious. Accidental line skipping has been referred to with a rather awkward French expression, *saut de ligne au hasard* ("jump from line to another"), in non French literature [28]. Because of the terms *homoeoteleuton* and *homoiotrope* have come to designate beginning and end of line similarities respectively, and by extension, accidental line skipping because of them [15].

This paper is organized as follows. Section 2 mentions some related work. Section 3 provides an outline of the KP algorithm's operation, necessary to understand how our extension works. Section 4 describes our extension, and Section 5 presents some experimental results.

#### 2 Related Work

Franklin [13] mentions the similarity problem in a survey of existing alternative TeX engines and remaining issues [13]. No solution is proposed in this paper, as it is merely a state-of-the-art review. Alex Hellner addresses the problem in his multiple-objective approach to line breaking [7]. However, the paper only seems to be concerned with beginning-of-line skipping (called *skewness*). Although interesting in theory, such an approach is not well suited to TeX. The paper does not provide a precise definition for stacks, and does not use the corresponding objective function in the reported experimental results. Other extensions to the KP algorithm have been proposed for the past, some with typographic [17], others with readability [18] in mind. For the latter, the main idea is to add a different layout [8]. Concerning the former, our underlying motivations are in fact opposite. The work in question attempts to provide flexibility at the expense of quality in order to cope with situations in which manual intervention is impossible, such as automatically adjusting to different displays. We, on the other hand, are interested

Permissions to make digital and hard copies of all or part of the work described or contained in this document in its entirety for private, internal use only, is granted by the author(s) for personal research purposes only, and is not to be resold in any form without the express written permission of the author(s). The work may not be copied, reproduced, or distributed in whole or in part by electronic means, or otherwise, without the prior permission of the author(s). Requests for permission should be addressed to Request permissions from permissions@acm.org. Downloaded 24 August 2024, from <https://doi.org/10.1145/3685408.3685666>. ACM holds the copyright for this work. Publication rights licensed to ACM. © 2024 Association for Computing Machinery. Publication rights licensed to ACM. ACM ID 1079-4807/24/0802-0001-0004 \$15.00. https://doi.org/10.1145/3685408.3685666



<https://didierverna.net/publications/verna.24.doceng/>



# Similarités

In olden times when wishing still helped one, there lived a king whose daughters were all beautiful; and the youngest was so beautiful that the sun itself, which has seen so much, was astonished whenever it shone in her face. Close by the king's castle lay a great dark forest, and under an old lime-tree in the forest was a well, and when the day was very warm, the king's child went out into the forest and sat down by the side of the cool fountain; and when she was bored she took a golden ball, and threw it up on high and caught it; and this ball was her favorite plaything.

▶ Solution de T<sub>E</sub>X

- ▶ Largeur : 201 pt
- ▶ Caractère : Latin Modern Roman 10 pt

▶ Problèmes (également en début de ligne)

- ▶ Micro-interruption
- ▶ Saut de ligne accidentel

▶ Vieille question (moines scribes)

- ▶ Cf. philologie/analyse de texte

▶ Terminologie

- ▶ « Saut du même au même » 😊
- ▶ « Cascade » (Thomas Savary)
- 💡 Homéoarchie/Homéotèleute (rhétorique)
- 💡 Échelle de caractères/mots



# Solutions

## 👊 Espaces insécables (cf. exposé de Thomas Savary « Composer comme un pro »)

- ▶ Homéoarchie: ~truc ... truc~
- ▶ Homéotélete: truc~ ... ~truc

## 💪 Un peu (plus) de douceur (dans ce mo...)

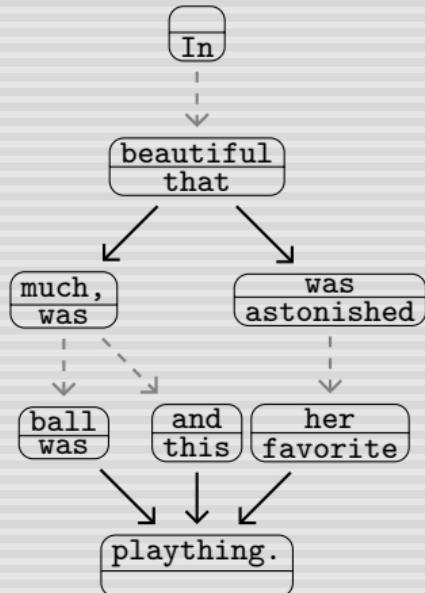
- ▶ Homéoarchie: \penaltyxxx truc ... truc\penaltyxxx
- ▶ Homéotélete: truc\penaltyxxx ... \penaltyxxx truc

## ▶ Dans tous les cas

- ▶ Travail manuel
- ▶  $xxx = ??$
- ▶ Non automatisable

## ▶ Automatiser !

# Le Knuth-Plass en 30 secondes



- ▶ Problème de type « plus court chemin »  
*Trouver la « meilleure » route du début jusqu'à la fin du paragraphe*
- ▶ Technique d'optimisation dynamique  
*Ne jamais construire l'intégralité du graphe*
- ▶ Fonction de coût
  - ▶ Démérites locaux : ligne par ligne  
*Badness, pénalité de césure, etc.*
  - ▶ Démérites contextuels : ligne contre ligne  
*Échelles de césure, disparités d'espacement, etc.*
- ▶ Décision finale :  $\min \sum(d. \text{ locaux} + \text{contextuels})$



# Extension au Knuth-Plass



Introduction



Knuth-Plass



Extension



Expérimentation



Conclusion

- ▶ **Idée :** un nouveau type de démerites contextuels (« démerites de similarité »)
  - ▶ Comparer les débuts & fins de lignes consécutives
  - ▶ Ajouter les démerites de similarité le cas échéant
  - ▶ Note : cas particulier des deux lignes finales
- ▶ **Implémentation :**
  - ▶ Sauvegarder les séquences de début/fin de ligne dans les noeuds
  - ▶ Seulement jusqu'à la première « glue » (espace élastique)/premier « dictionnaire » (point de césure/ligature)
  - ▶ Supprimer les « kerns »
- ▶ **Justification :**
  - ▶ Dictionnaires : éviter le surcoût de la déconstruction 🤔
  - ▶ Glue : éviter les considérations d'alignement vertical
  - ▶ Kerns : petits ajustements, identiques en cas de similarité



# Extension au Knuth-Plass



Introduction



Knuth-Plass



Extension



Expérimentation



Conclusion

▶ **Idée :** un nouveau type de démerites contextuels (« démerites de similarité »)

- ▶ Comparer les débuts & fins de lignes consécutives
- ▶ Ajouter les démerites de similarité le cas échéant
- ▶ Note : cas particulier des deux lignes finales

▶ **Implémentation :**

- ▶ Sauvegarder les séquences de début/fin de ligne dans les noeuds
- ▶ Seu' **Comparer des séquences courtes**
- ▶ Sup. *Qui peut le moins peut le plus...*

naire » (point de c

▶ **Justification :**

- ▶ Discrétionnaires : éviter le surcoût de la déconstruction 🤔
- ▶ Glue : éviter les considérations d'alignement vertical
- ▶ Kerns : petits ajustements, identiques en cas de similarité



# Démonstration



Introduction



Knuth-Plass



Extension



Expérimentation



Conclusion

In olden times when wishing still helped one, there lived a king whose daughters were all beautiful; and the youngest was so beautiful that the sun itself, which has seen so much, was astonished whenever it shone in her face. Close by the king's castle lay a great dark forest, and under an old lime-tree in the forest was a well, and when the day was very warm, the king's child went out into the forest and sat down by the side of the cool fountain; and when she was bored she took a golden ball, and threw it up on high and caught it; and this ball was her favorite plaything.

\similar{demerits=0}

In olden times when wishing still helped one, there lived a king whose daughters were all beautiful; and the youngest was so beautiful that the sun itself, which has seen so much, was astonished whenever it shone in her face. Close by the king's castle lay a great dark forest, and under an old lime-tree in the forest was a well, and when the day was very warm, the king's child went out into the forest and sat down by the side of the cool fountain; and when she was bored she took a golden ball, and threw it up on high and caught it; and this ball was her favorite plaything.

\similar{demerits=2800}

In olden times when wishing still helped one, there lived a king whose daughters were all beautiful; and the youngest was so beautiful that the sun itself, which has seen so much, was astonished whenever it shone in her face. Close by the king's castle lay a great dark forest, and under an old lime-tree in the forest was a well, and when the day was very warm, the king's child went out into the forest and sat down by the side of the cool fountain; and when she was bored she took a golden ball, and threw it up on high and caught it; and this ball was her favorite plaything.

\similar{demerits=5230}



# Valeur ajoutée



Introduction



Knuth-Plass



Extension



Expérimentation



Conclusion

- ▶ **Pertinence** : le problème est-il fréquent ?
- ▶ **Efficacité** : cette extension le résout-il ?
- ▶ **Deux expériences orthogonales :**
  - ▶ Un paragraphe à des largeurs très variées
  - ▶ De nombreux paragraphes à une largeur fixe

## Conditions expérimentales

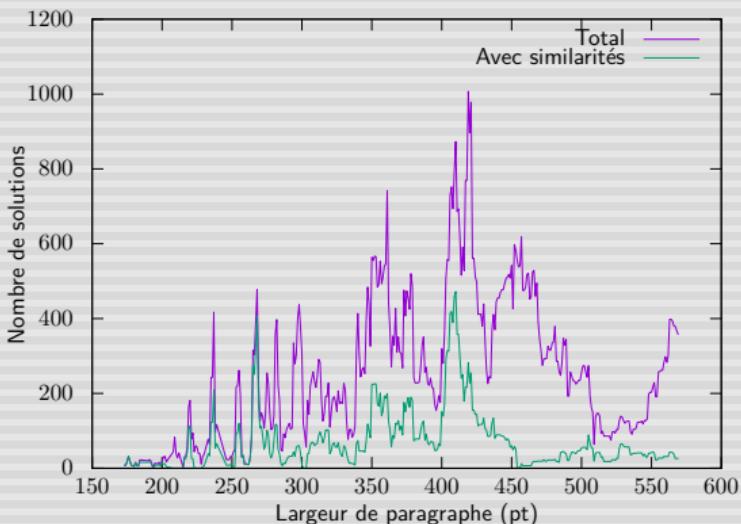
- ▶ *Frog King* (des Frères Grimm) §1
- ▶ 142 pt ( $\approx 5$  cm)  $\rightarrow$  569 pt ( $\approx 20$  cm)
- ▶ 427 cas
- ▶ Total :
  - ▶ Expériences 1 & 2  $\rightarrow$  1 951 cas
  - ▶  $\times 3 \rightarrow 5 853$  passes

## Conditions expérimentales

- ▶ *Moby Dick* (Herman Melville)
- ▶ 1 524 paragraphes à 284 pt ( $\approx 10$  cm)

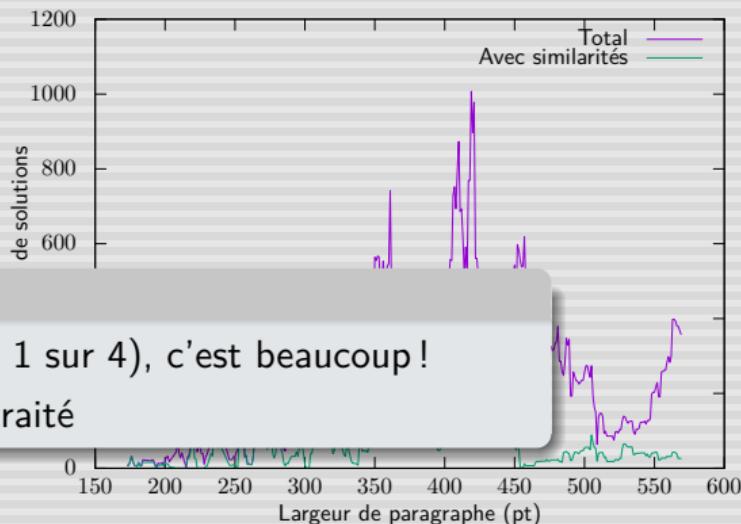
# Pertinence (*Frog King*)

- ▶ Résultats attendus
  - ▶ Formes chaotiques
  - ▶ Valeurs proches pour des paragraphes étroits
- ▶ Découvertes
  - ▶ Similarités fréquentes
  - ▶ Solutions alternatives également
- ▶ Comportement de T<sub>E</sub>X
  - ▶ Similarités inévitables : 4 %
  - ▶ Solution contenant des similarités
    - ▶ Expérience 1 : 21 %
    - ▶ Expérience 2 : 26 %



# Pertinence (*Frog King*)

- ▶ Résultats attendus
  - ▶ Formes chaotiques
  - ▶ Valeurs proches pour des paragraphes étroits
- ▶ Découvertes
  - ▶ Sim
  - ▶ Solu
- ▶ Comport
  - ▶ Sim
  - ▶ Solution contenant des similarités
    - ▶ Expérience 1 : 21 %
    - ▶ Expérience 2 : 26 %



## Conclusion

# Efficacité

$\backslash similardemerits=10000$	$+ \backslash adjacentdemerits=0$
Corrigés : 48 %/50 % Améliorés : 50 %/63 %	Corrigés : 53 %/66 % Améliorés : 57 %/73 %

- ▶ Conclusion : Traitement automatique possible

! Attention !

- ▶ Colonne de droite = conditions extrêmes
- ▶ Moins de similarités ne signifie pas forcément meilleure esthétique
- ▶ Mais il y a une marge de manœuvre



# Conclusion

- ▶ Traitement des similarités pour le Knuth-Plass
  - ▶ Extension simple & légère
  - ▶ Rétrocompatible avec  $\text{\TeX}$  (`\similar{demerits}=0`)
- ▶ Implémenté dans ETAP, utile dans les moteurs en production
- ▶ Expérimentation → traitement automatique des similarités possible et souhaitable

## Perspectives

- ▶ Étude de l'interaction avec les autres critères esthétiques
- ▶ Démérites de similarité en tant que fonction plutôt que valeur scalaire
- ▶ LuaMeta $\text{\TeX}$  — `\left/right twindemerits`