# CASIMIR : un Corpus d'Articles Scientifiques Intégrant les ModIfications et Révisions des auteurs

Léane Jourdan, Florian Boudin, Richard Dufour, Nicolas Hernandez
{prénom.nom}@univ-nantes.fr

## Objectifs

- Nouveau corpus d'articles scientifiques avec révisions
- Un alignement des modifications entre les versions des articles
- Enrichit avec les métadonnées des articles et les relectures par les pairs

## Processus de collecte

1. Collecte des événements (ateliers, conférences, etc) sur Open Review
2. Collecte des articles : métadonnées, relectures et PDF disponibles des versions des articles
3. Filtrage des articles ayant une seule version (89,33% des articles conservés soit 97,46% des PDF)
4. Conversion des PDF vers XML (outil Grobid)

## Exemple de révisions



## Contenu

### Articles
les PDF des versions des articles

### Fichiers de correspondance
- entre les versions finales et antérieures des articles
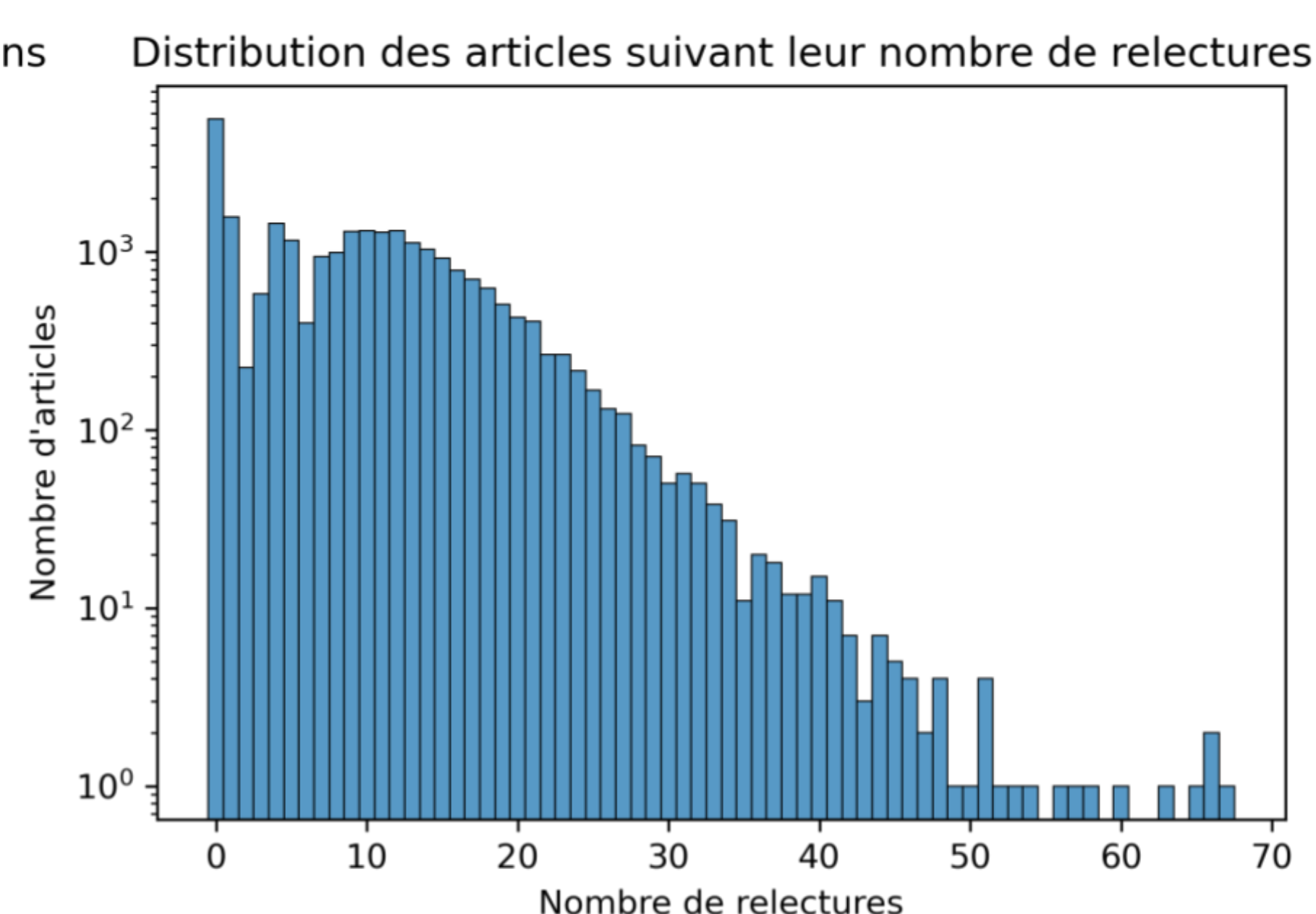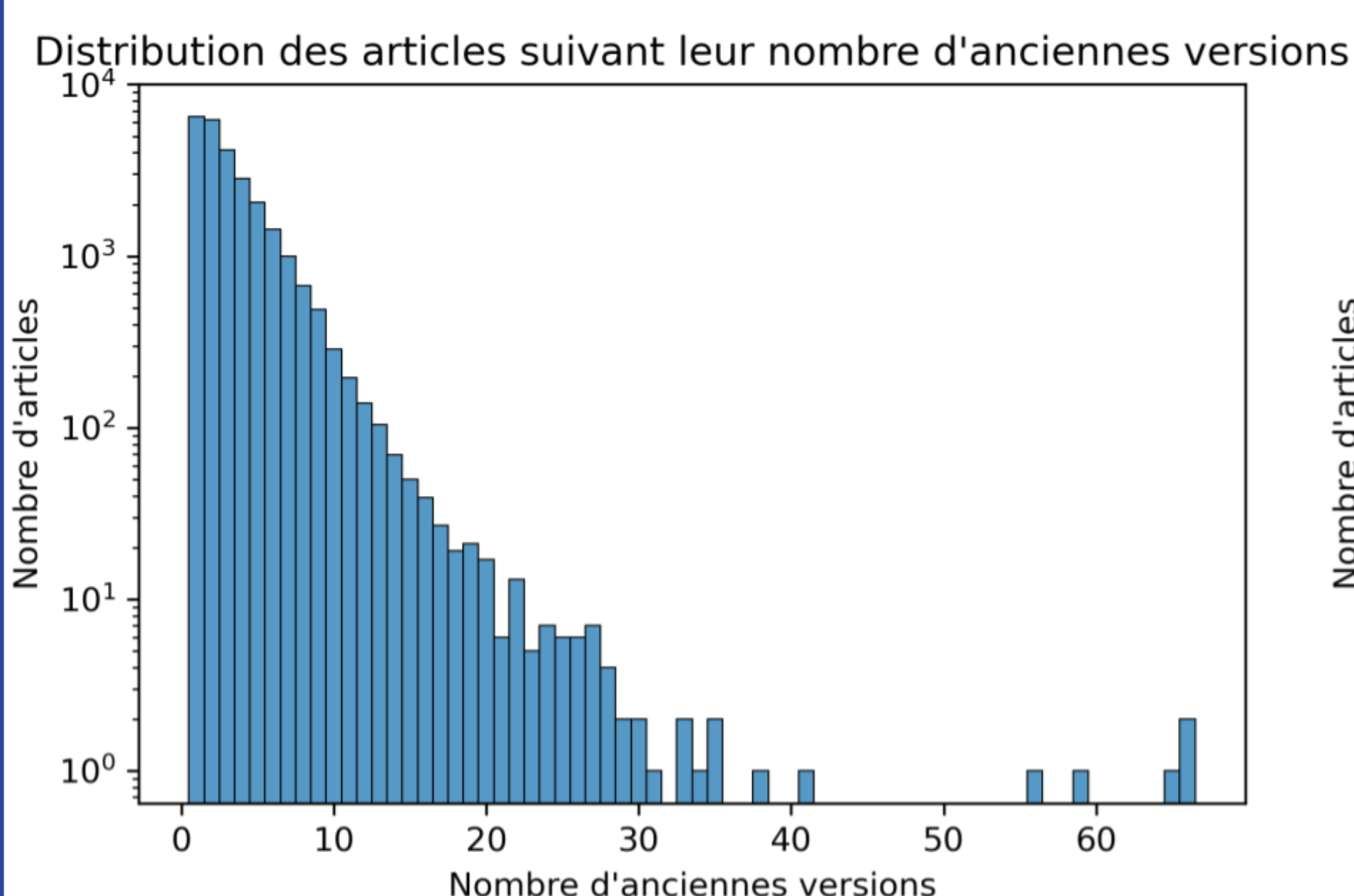- entre les relectures et les articles

### Metadonnées des articles
- dates
- auteurs
- mots-clés
- événement
- identifiants

### Relectures
- commentaires
- notes
- décisions
- dates

## Description

Contient 730 évènements et 118 415 PDF pour 26 355 articles
Domaines : apprentissage auto, robotique, TAL, vision, etc.



Distribution des articles suivant leur nombre d'anciennes versions



Distribution des articles suivant leur nombre de relectures

## Et ensuite?

- Améliorer la conversion des PDF (choix de l'outil, traitement des figures, tables)
- Aligner les versions paragraphe à paragraphe et phrase à phrase puis extraire les révisions
- Annoter les documents selon une taxonomie de révisions à définir (ex: clarté, grammaire, style)
- Exploitation pour la mise en place d'outils d'aide à l'écriture