



Orphan Works As Grist For The Data Mill

Matthew Sag

Associate Professor, Loyola University Chicago School of Law

Paper available available at <http://ssrn.com/abstract=2038889>).

Three Faces of Library Digitization

◆ Preservation

◆ Data production and analysis

- Searching books, testing search algorithms, computational linguistics, automated translation, natural language processing, macro-analysis of text

◆ A platform for **display and distribution** of individual works

Library digitization and orphan works

◆ Key Question:

- Does copying for a ~~non-consumptive~~ **nonexpressive** use implicate the rights of the copyright owner?

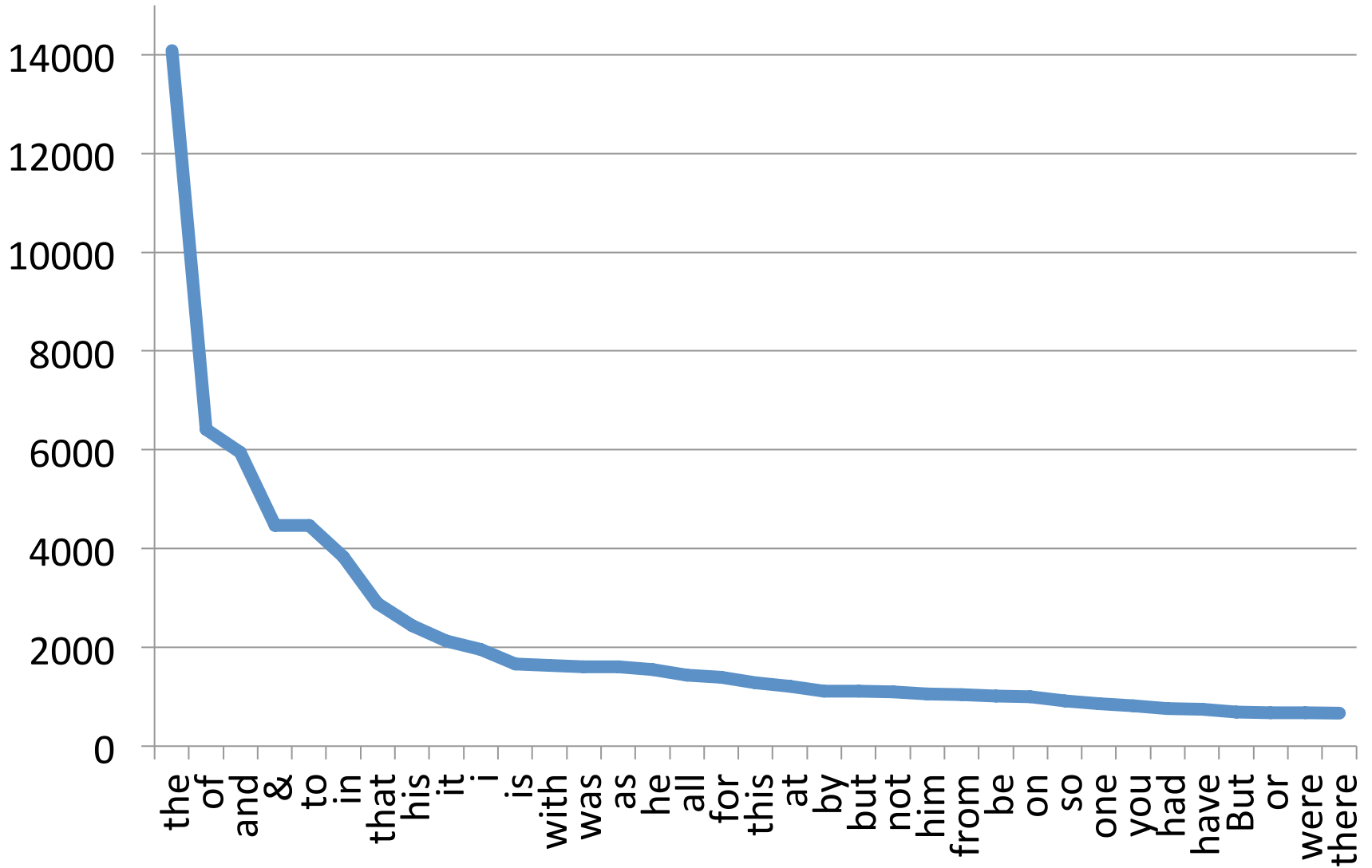
◆ Note:

- Orphan works explains why we care, but the orphan status of these works is not directly relevant to the primary question.

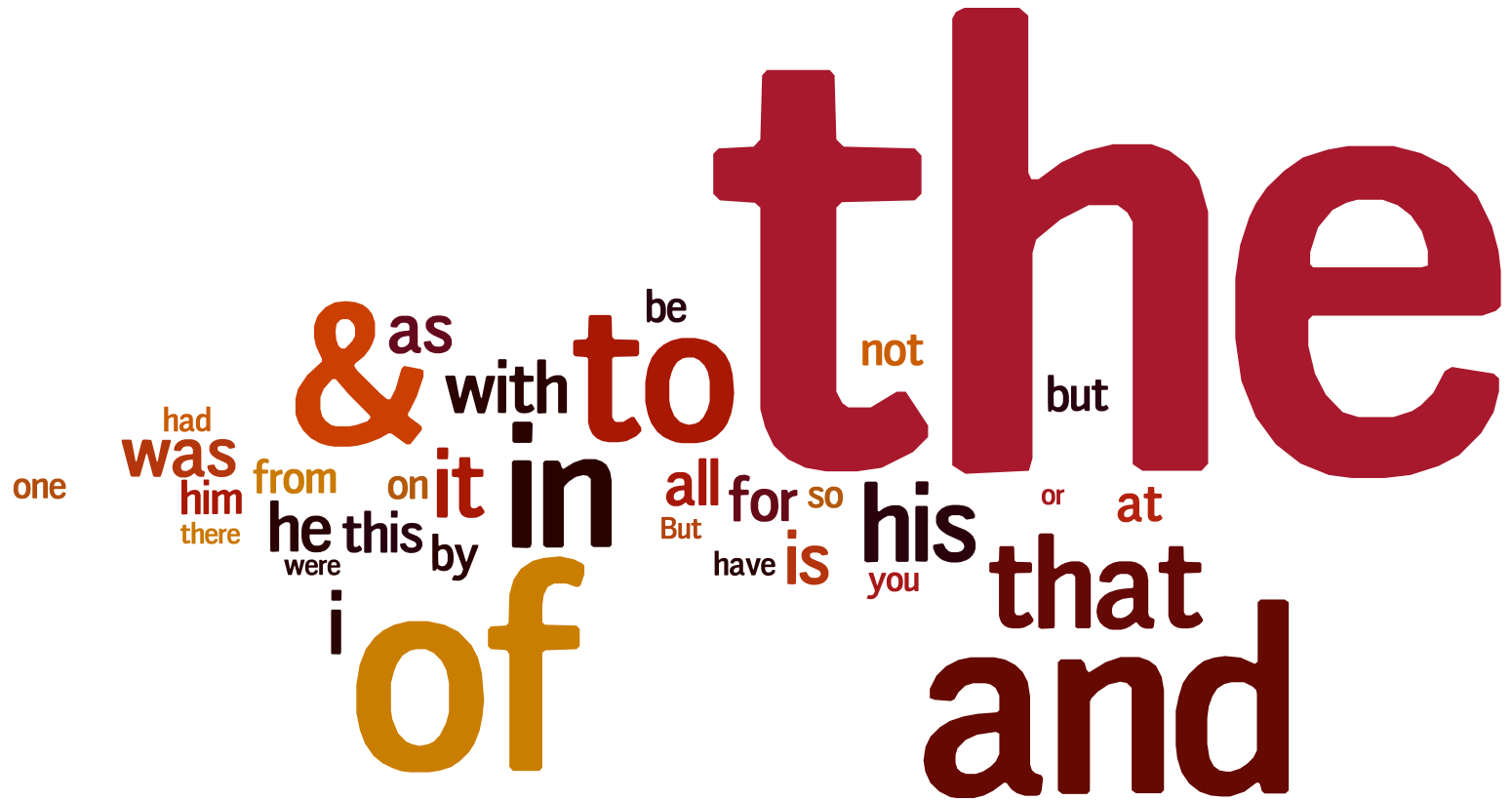
Thought Experiment

- ◆ Brian is a savant with total recall
- ◆ Moby Dick has its copyright restored
 - (Perpetual Copyright Act of 2014??)
- ◆ Brian produces a frequency table

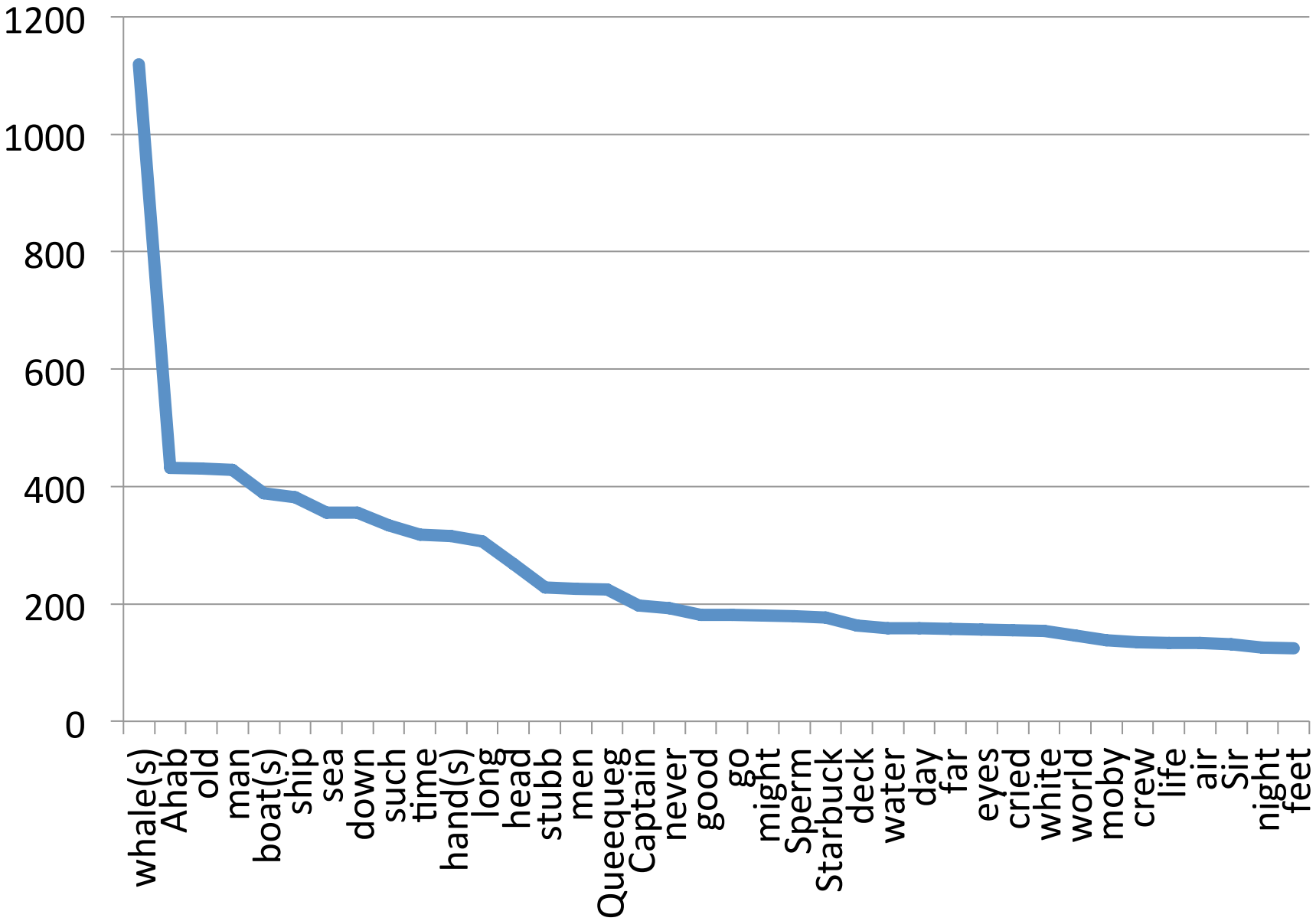
Common words in Moby Dick



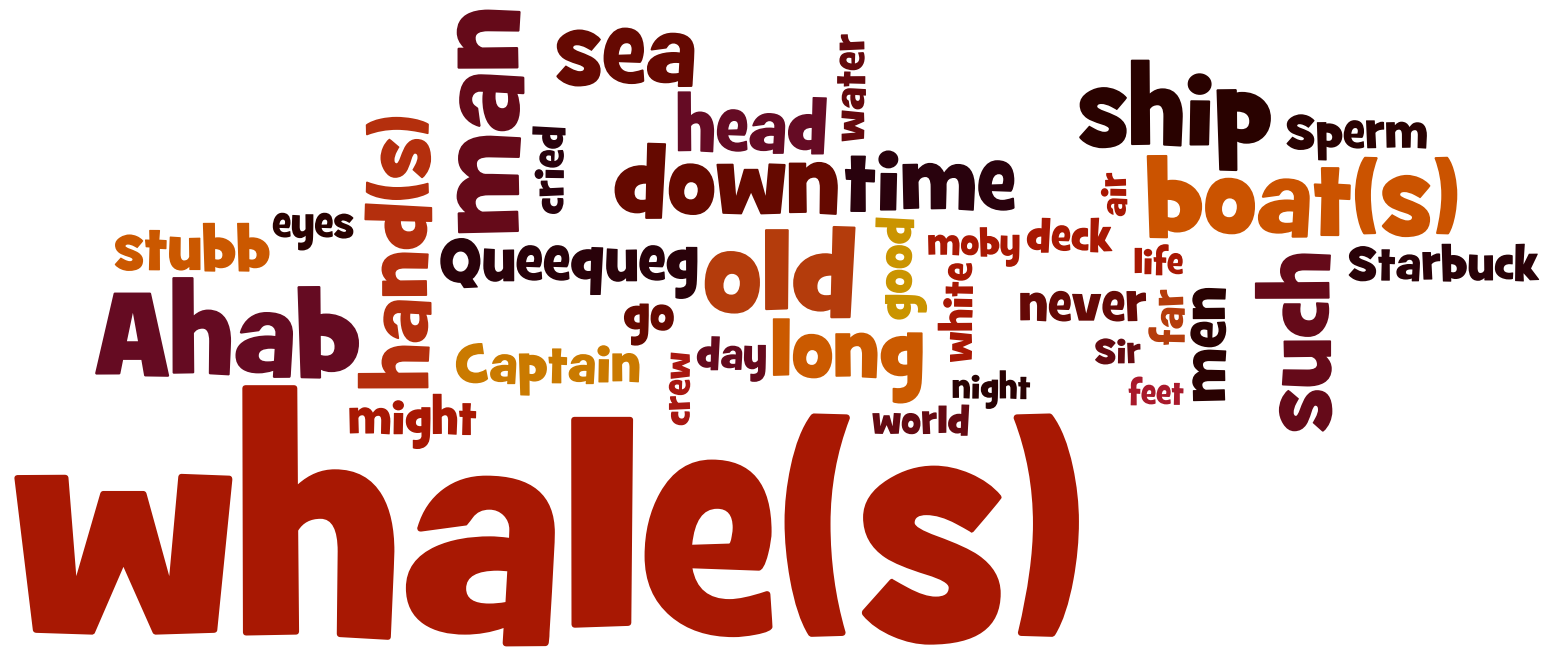
Common words in Moby Dick



Uncommon words in Moby Dick



Uncommon words in Moby Dick



Meta Data – a restatement of the obvious

- ◆ **Meta data (even if its valuable) does not infringe the rights of the copyright owner.**
 - Idea-expression distinction
 - Merger
 - Substantially similarity →
 - Originality → →

Originality

- [1] **"Goblin-made armour does not require cleaning, simple girl. Goblins' silver repels mundane dirt, imbibing only that which strengthens it."** (J.K. Rowling, Deathly Hallows)
- [2] **".. goblin-made armor does not require cleaning, because goblins' silver repels mundane dirt, imbibing only that which strengthens it, such as basilisk venom."** (Harry Potter Lexicon)
- [3] **Other than 'Goblin', none of the words in [1] are repeated.** (Matthew Sag)
- [4] **There is a high level of similarity between [1] and [2](anti-plagiarism software)**

Producing Meta Data – Not quite so obvious

- ◆ Hard to argue that a reading machine (e.g. Google Book Search) does not ‘reproduce the work’ in a ‘copy’, even if no one reads it.
- ◆ The distinction between expressive and nonexpressive works is well recognized. The same distinction should generally be made in relation to potential acts of infringement.
 - Copying for purely nonexpressive purposes, such as the automated extraction of data, should not be regarded as infringing.

Statutory rights of the author are limited to the communication of original expression to the public

◆ Consider

- Threshold of substantial similarity is defined in reference to the perspective of the ordinary observer (with some filtering of facts, ideas, etc.).
- Intermediate copying does not infringe (screen-play cases), is fair use (reverse engineering cases)

Implications

- ◆ **Automated reproduction for nonexpressive uses** (such as search engines, plagiarism detection, and macro-literary analysis) **does not communicate the author's original expression to the public**
 - No expressive substitution, no infringement

Application to Fair Use

- ◆ **(1) purpose and character:** Like transformative uses, a nonexpressive use poses no risk of expressive substitution
- ◆ (2) nature of the work ... “not much use”
- ◆ **(3) Amount and Substantiality:** Like transformative uses, because there is no expressive substitution in a nonexpressive use, the amount of copying is qualitatively insignificant.
- ◆ **(4) Market effect:** Like transformative uses, a nonexpressive use poses no risk of expressive substitution, thus no cognizable market effect.

