

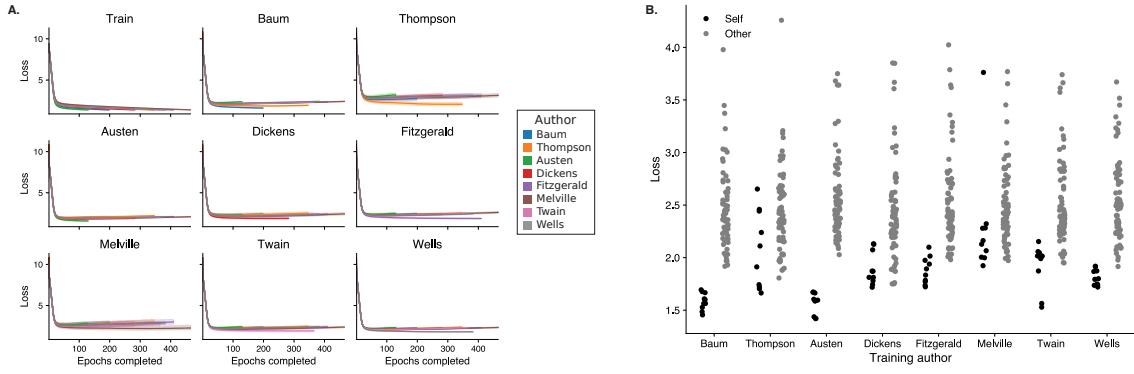
*Supplementary materials for: A Stylometric Application of
Large Language Models*

Harrison F. Stropkay, Jiayi Chen, Daniel N. Rockmore, and Jeremy R. Manning

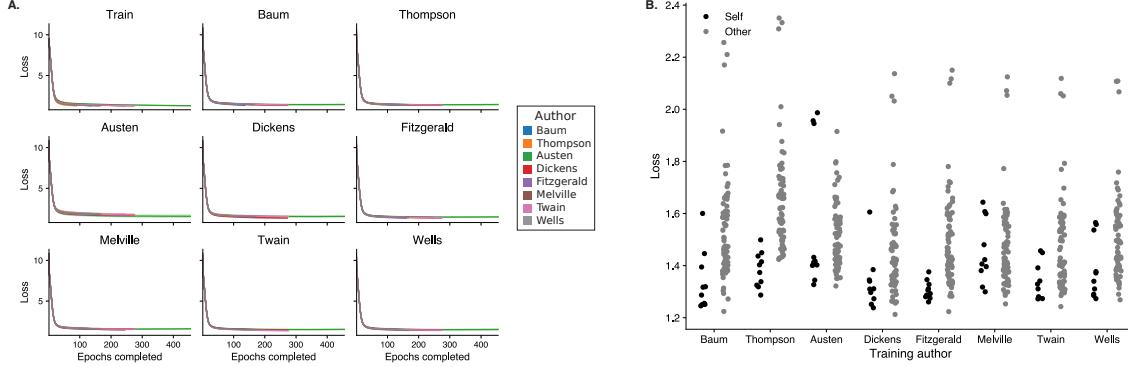
Dartmouth College

Hanover, NH 03755, USA

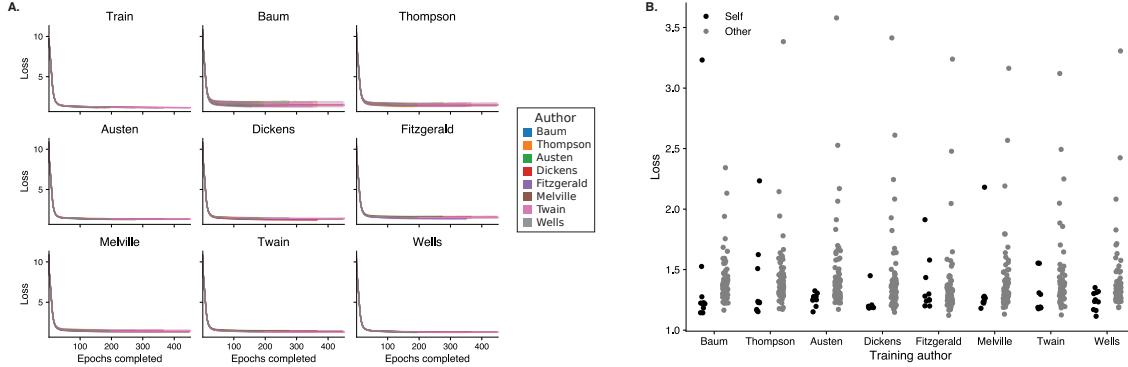
{harrison.f.stropkay.25, jiayi.chen.gr,
daniel.n.rockmore, jeremy.r.manning}@dartmouth.edu



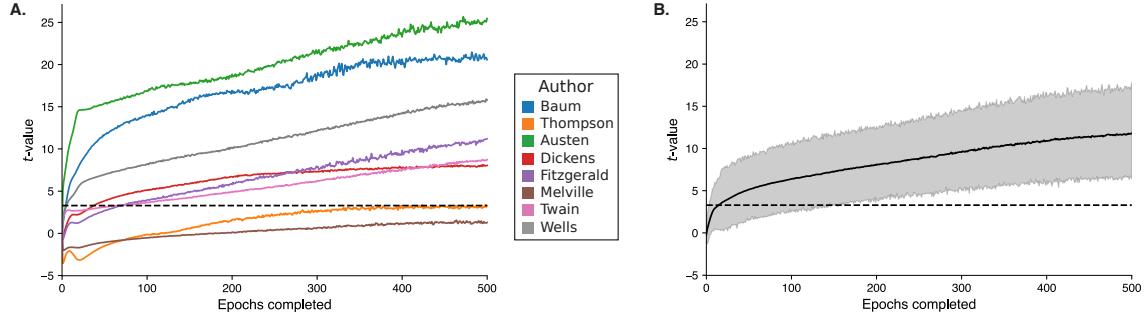
Supplementary Figure 1: Cross-entropy loss across models and authors using only content words. Follows the general format of Figure 1 in the main text, but uses models trained on only content words. All function words are masked out using <FUNC>. **A.** Average cross-entropy loss on *Training* data and held-out test data from each author, plotted as a function of the number of training epochs. Each color denotes a model trained on a single author’s work. Error ribbons denote bootstrap-estimated 95% confidence intervals over 10 random seeds. **B.** Cross-entropy loss assigned to held-out test data by each author’s model (*x*-axis). Held-out test data is either from the *same* author (black) or from *other* authors (gray). Each dot denotes the average loss (across all 1024-token chunks) for a single random seed.



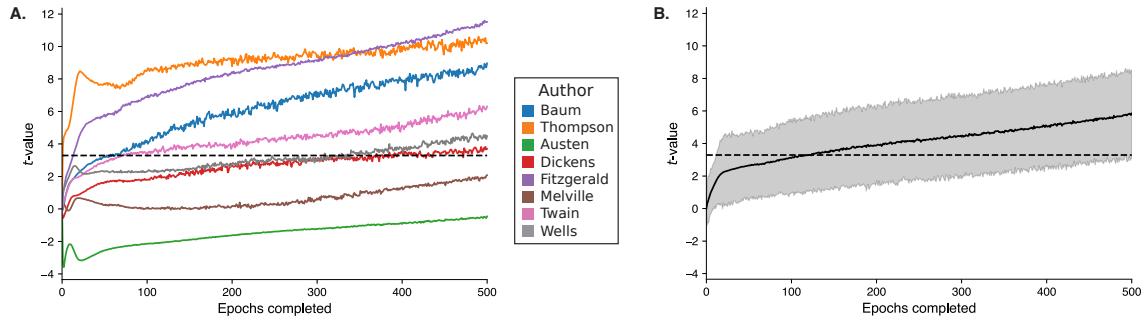
Supplementary Figure 2: Cross-entropy loss across models and authors using only function words. Follows the general format of Figure 1 in the main text, but uses models trained on only function words. All content words are masked out using <CONTENT>. **A.** Average cross-entropy loss on *Training* data and held-out test data from each author, plotted as a function of the number of training epochs. Each color denotes a model trained on a single author's work. Error ribbons denote bootstrap-estimated 95% confidence intervals over 10 random seeds. **B.** Cross-entropy loss assigned to held-out test data by each author's model (*x*-axis). Held-out test data is either from the *same* author (black) or from *other* authors (gray). Each dot denotes the average loss (across all 1024-token chunks) for a single random seed.



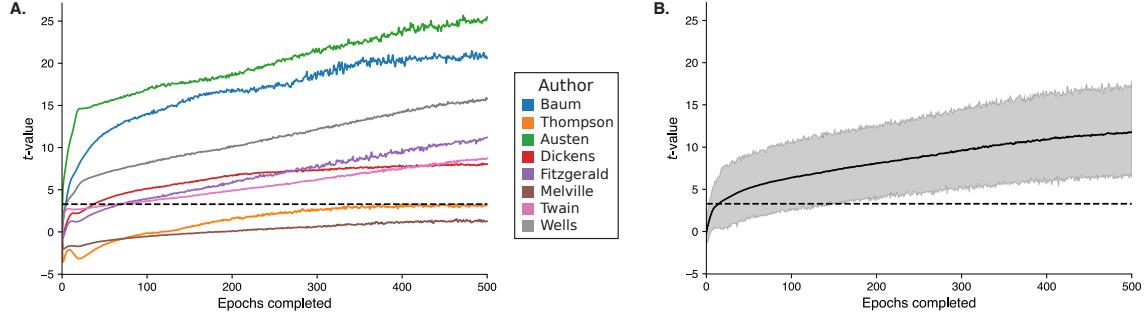
Supplementary Figure 3: Cross-entropy loss across models and authors using only parts of speech. Follows the general format of Figure 1 in the main text, but uses models trained on only parts of speech. All words are replaced with their corresponding part of speech tag. **A.** Average cross-entropy loss on *Training* data and held-out test data from each author, plotted as a function of the number of training epochs. Each color denotes a model trained on a single author's work. Error ribbons denote bootstrap-estimated 95% confidence intervals over 10 random seeds. **B.** Cross-entropy loss assigned to held-out test data by each author's model (*x*-axis). Held-out test data is either from the *same* author (black) or from *other* authors (gray). Each dot denotes the average loss (across all 1024-token chunks) for a single random seed.



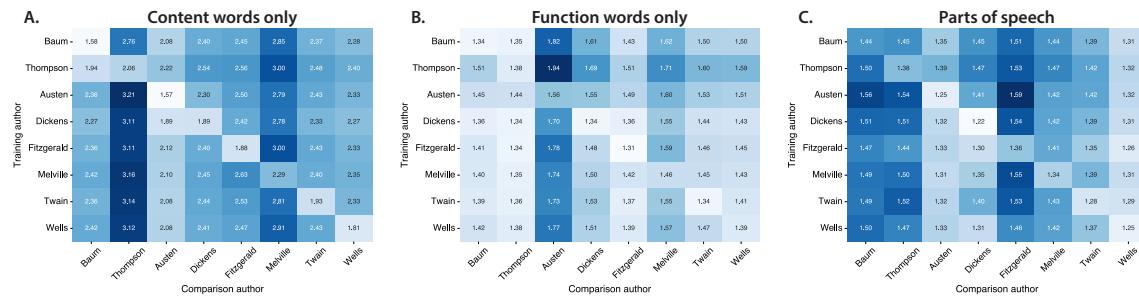
Supplementary Figure 4: Same vs. other author comparisons, by model, using only content words. Follows the general format of Figure 2in the main text, but uses models trained on only content words. All function words are masked out using <FUNC>. **A.** Each curve denotes, as a function of the number of training epochs, the the t -statistic from a t -test comparing the distribution of losses (across random seeds) assigned to held-out texts from the given author (color) versus held-out texts from all other authors. **B.** The average t -statistic across all eight authors, as a function of the number of training epochs. Error ribbons denote bootstrap-estimated 95% confidence intervals across authors.



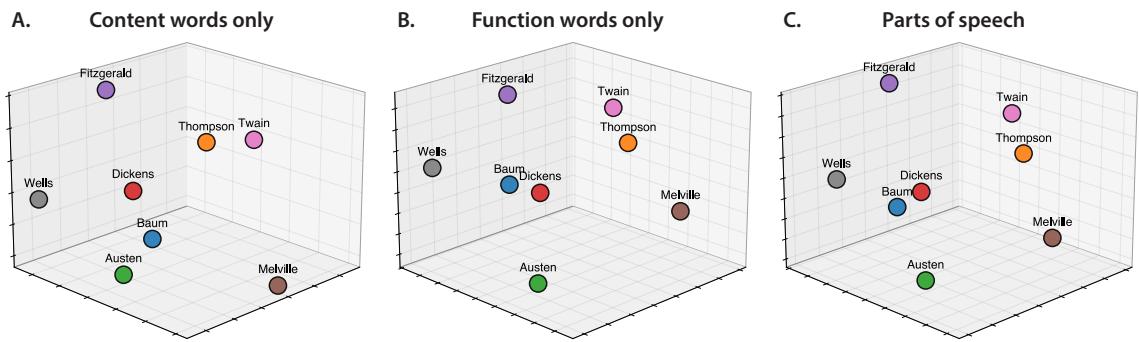
Supplementary Figure 5: Same vs. other author comparisons, by model, using only function words. Follows the general format of Figure 2in the main text, but uses models trained on only function words. All content words are masked out using <CONTENT>. **A.** Each curve denotes, as a function of the number of training epochs, the the t -statistic from a t -test comparing the distribution of losses (across random seeds) assigned to held-out texts from the given author (color) versus held-out texts from all other authors. **B.** The average t -statistic across all eight authors, as a function of the number of training epochs. Error ribbons denote bootstrap-estimated 95% confidence intervals across authors.



Supplementary Figure 6: Same vs. other author comparisons, by model, using only parts of speech. Follows the general format of Figure 2in the main text, but uses models trained on only parts of speech. All words are replaced with their corresponding part of speech tag. **A.** Each curve denotes, as a function of the number of training epochs, the the t -statistic from a t -test comparing the distribution of losses (across random seeds) assigned to held-out texts from the given author (color) versus held-out texts from all other authors. **B.** The average t -statistic across all eight authors, as a function of the number of training epochs. Error ribbons denote bootstrap-estimated 95% confidence intervals across authors.



Supplementary Figure 7: Confusion matrices. Follows the general format of Figure 3in the main text, but shows confusion matrices for models trained on only content words (A), only function words (B), and only parts of speech (C). Within each panel, the matrix displays the average cross-entropy loss assigned by models trained on each author's writing (column) to held-out texts from each author (row), after subtracting the native author's baseline loss.



Supplementary Figure 8: Multidimensional scaling plot. Follows the general format of Figure 4 in the main text, but shows MDS projections of the (symmetrized) average cross entropy loss matrices shown in Figure 7, for models trained on only content words (A), only function words (B), and only parts of speech (C).