# SUPPLEMENTS
## Mode Classification and Natural Units in Plainchant
### Bas Cornelissen, Willem Zuidema and John Ashley Burgoyne

**Data and code.** All data and code used in this study has been made available online at github.com/bacor/ISMIR2020. All randomness in the code has been fixed, so it should in theory be possible to reproduce our results exactly. The evaluations metrics of all experiments are already included in the repository, as is the data used in the first run of the experiment; this should be sufficient for reproducing most figures. We have included model predictions and tuning results only for the first run of the experiment. Detailed logs of everything from data generation to vizualization can also be in the repository, together with many more figures besides those included in the paper and the supplements. In particular, the repository contains heatmaps with multiple evaluation metrics (accuracy, precision, recall and $F_1$) for all models and all experimental conditions.

```
Filtering chants...
. Filter Chants Without Volpiano:
. Exclude all chants with an empty volpiano field
.  > 87.20% removed (433443 out of 497071; 63628 remain)
. Filter Chants Without Notes:
. Exclude all chants without notes
.  > 2.87% removed (1825 out of 63628; 61803 remain)
. Filter Chants Without Simple Mode:
. Include only chants with simple modes: 1-8, not transposed
.  * include_transposed=False
.  > 23.02% removed (14227 out of 61803; 47576 remain)
. Filter Chants Without Full Text:
. Filter chants without full text
.  > 20.65% removed (9823 out of 47576; 37753 remain)
. Filter Chants Where Incipit Is Full Text:
. Filter chants whose incipit is identical to the full text
.  > 14.59% removed (5507 out of 37753; 32246 remain)
. Filter Chants By Genre:
. Include only chants with a certain genre
.  * include=['genre_a']
.  > 52.06% removed (16787 out of 32246; 15459 remain)
. Filter Chants Not Starting With G Clef:
. Exclude chants that do not start with a G clef
.  > 0.05% removed (7 out of 15459; 15452 remain)
. Filter Chants With F Clef:
. Exclude chants that contain an F clef
.  > 0.00% removed (0 out of 15452; 15452 remain)
. Filter Chants With Missing Pitches:
. Filter chants with missing pitches: containing the substring 6------6
.  > 7.54% removed (1165 out of 15452; 14287 remain)
. Filter Chants With Nonvolpiano Chars:
. Exclude all chants with non-volpiano characters
.  > 0.03% removed (5 out of 14287; 14282 remain)
. Filter Chants Without Word Boundary:
. Only include chants with '---' in their volpiano
.  > 0.08% removed (11 out of 14282; 14271 remain)
. Filter Chants With Duplicated Notes:
. Filter duplicate chants: whose notes occur multiple times
.  > 2.84% removed (406 out of 14271; 13865 remain)
```
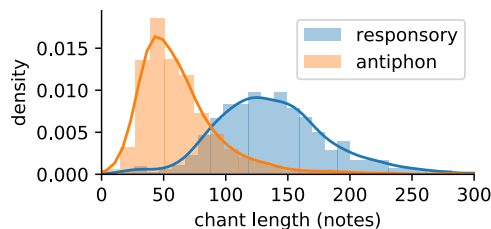
**Figure S1.** **Filtering.** As described in the main text, we filtered the total dataset of 497,071 chants to obtain a clean subset of responsories and antiphons. The effects of all of the filters are logged and will be made available online. As an example, this 'figure' shows the series of filters applied to obtain the full set of antiphons used in this study.

| Genre | Subset | Split | # chants | # notes | Mean length (notes) |
|---|---|---|---|---|---|
| responsory | full | train | 4 922 | 676 807 | 137.5 |
| responsory | full | test | 2 109 | 290 064 | 137.5 |
| responsory | full | total | 7 031 | 966 871 | 137.5 |
| responsory | subset | train | 1 234 | 169 642 | 137.5 |
| responsory | subset | test | 529 | 72 504 | 137.1 |
| responsory | subset | total | 1 763 | 242 146 | 137.3 |
| antiphon | full | train | 9 706 | 576 738 | 59.4 |
| antiphon | full | test | 4 159 | 248 405 | 59.7 |
| antiphon | full | total | 13 865 | 825 143 | 59.5 |
| antiphon | subset | train | 2 911 | 190 165 | 65.3 |
| antiphon | subset | test | 1 248 | 82 781 | 66.3 |
| antiphon | subset | total | 4 159 | 272 946 | 65.6 |

**Figure S2**. **Dataset statistics.** The number of chants, their average length and number of notes for each dataset. We sort datasets by genre, then by subset (include melody variants in the full set, or exclude them in the subset) and finally by train/test split (or total for the two combined). The train/test splits are different in each run of the experiment. These statistics are computed from the data used in the first run, and others are comparable.

| genre | dataset | kind | top mode | frequency |
|---|---|---|---|---|
| responsory | full | train | 8 | 20.85% |
| responsory | full | test | 8 | **21.13**% |
| responsory | subset | train | 1 | 21.65% |
| responsory | subset | test | 1 | 20.19% |
| antiphon | full | train | 8 | 28.47% |
| antiphon | full | test | 8 | **28.13**% |
| antiphon | subset | train | 1 | 23.50% |
| antiphon | subset | test | 1 | 24.18% |

**Figure S3**. **Majority baselines**. The frequency of the largest classes in each of the datasets. Boldfaced values correspond to the classification *accuracy* of the worst-performing conditions discussed in the main text. (The frequencies are marginally different in the five experimental runs; shown are the averages.)



**Figure S4**. **Chant lengths in two genres.** Responsories are usually much longer than antiphons. The distribution is estimated from the training datasets without melody variants.

**Figure S5**. **Lengths of natural units.** Natural units have different lenghts in responsories and antiphons. Responsories are more *melismatic*: they use more notes per syllable. As a result, a typical word is also much longer. This is shown in the figure using violin plots, a visualization of the length distribution using a kernel density estimate. Note that the total area has no meaning in this plot; we normalized the widths of the violins for better readability. The distributions are estimated from the training datasets without melody variants).

|  | neume | syllable | word |
|---|---|---|---|
| antiphon | 1.50 | 1.55 | 3.98 |
| responsory | 2.32 | 2.96 | 7.12 |

**Figure S6**. **Mean lengths of natural units.** Natural units have different lengths in responsories and antiphons, as the mean lenghts in number of notes shows. Figure S5 shows the full distribution. Means estimated from the training datasets without melody variants.



**Figure S7**. **Pitch class profiles** The pitch class profiles used in the profile approach (cf [12]). Shown are data for responsories, estimated from the training data without melody variants.

**Figure S8. Repetition profiles** The repetition profiles used in the profile approach. Every bar shows the average number of repetitions of that note in a chant (see main text for details). Shown are data for responsories, estimated from the training data without melody variants.

**Figure S9**. PCA **plots.** Two-dimensional representation of the high-dimensional feature space in several different conditions. A sample of TF–IDF vectors is shown, after reducing their dimensionality using a principal component projection.

**Figure S10**. *t*-SNE **plots.** Two-dimensional representation of the high-dimensional feature space in several different conditions. A sample of TF–IDF vectors is shown, after reducing their dimensionality using *t*-SNE, a nonlinear dimensionality-reduction technique that maximizes the probability of mapping nearby points in high-dimensional space to nearby points in a lower dimensional space. The axes have no natural interpretation, but the graphs suggest that clusters are most clearly separated at the left, and mostly overlapping at the right.

**Figure S11**. **Melody variants in Cantus.** The top three panels show examples of sets of melody variants: the first 100 notes of melodies sharing a Cantus ID. Different colors correspond to different pitches, or more precisely, different Volpiano characters after discarding dashes. As a comparison, the bottom panel shows 100 notes of 20 random melodies.

**A Classical approach**

| | responsory | antiphon |
|---|---|---|
| final | $39.7^{\pm1.2}$ | $48.6^{\pm0.8}$ |
| ambitus | $55.7^{\pm0.6}$ | $60.8^{\pm1.0}$ |
| initial | $37.5^{\pm0.4}$ | $47.8^{\pm1.5}$ |
| final ambitus | $88.8^{\pm0.5}$ | $79.5^{\pm0.7}$ |
| final initial | $73.3^{\pm1.0}$ | $72.1^{\pm0.5}$ |
| ambitus initial | $70.3^{\pm0.4}$ | $73.1^{\pm0.4}$ |
| final ambitus initial | $\underline{\mathbf{89.8}^{\pm0.6}}$ | $\underline{\mathbf{86.3}^{\pm0.6}}$ |

**B Profile approach**

| | responsory | antiphon |
|---|---|---|
| pitch class profile | $85.1^{\pm0.8}$ | $88.3^{\pm0.3}$ |
| pitch profile | $\underline{\mathbf{87.8}^{\pm0.7}}$ | $\underline{\mathbf{89.6}^{\pm0.2}}$ |
| repetition profile | $80.9^{\pm1.4}$ | $84.2^{\pm0.2}$ |

**C Distributional approach: responsories**

| | pitch | dep. interval | indep. interval | dep. contour | indep. contour |
|---|---|---|---|---|---|
| neumes | $92.1^{\pm0.4}$ | $86.2^{\pm0.6}$ | $79.0^{\pm0.3}$ | $62.9^{\pm1.1}$ | $52.2^{\pm0.7}$ |
| syllables | $\underline{\mathbf{93.1}^{\pm0.7}}$ | $\underline{\mathbf{88.6}^{\pm0.8}}$ | $85.8^{\pm0.6}$ | $78.8^{\pm0.9}$ | $76.2^{\pm2.1}$ |
| words | $90.4^{\pm1.0}$ | $86.9^{\pm0.5}$ | $\underline{\mathbf{86.2}^{\pm0.8}}$ | $\underline{\mathbf{82.3}^{\pm0.6}}$ | $\underline{\mathbf{80.9}^{\pm1.5}}$ |
| 1-mer | $86.6^{\pm0.4}$ | $53.4^{\pm0.8}$ | $7.4^{\pm0.4}$ | $20.0^{\pm0.6}$ | $7.4^{\pm0.4}$ |
| 2-mer | $90.5^{\pm0.5}$ | $74.0^{\pm0.7}$ | $37.6^{\pm1.2}$ | $25.2^{\pm0.5}$ | $17.1^{\pm0.3}$ |
| 3-mer | $91.8^{\pm0.6}$ | $80.5^{\pm0.5}$ | $65.4^{\pm0.9}$ | $36.7^{\pm0.5}$ | $22.7^{\pm0.5}$ |
| 4-mer | $91.5^{\pm0.4}$ | $83.2^{\pm1.0}$ | $75.3^{\pm1.3}$ | $47.2^{\pm0.5}$ | $34.0^{\pm0.9}$ |
| 5-mer | $90.5^{\pm1.0}$ | $83.6^{\pm0.8}$ | $80.5^{\pm0.5}$ | $53.5^{\pm0.8}$ | $42.6^{\pm1.1}$ |
| 6-mer | $88.0^{\pm1.1}$ | $82.5^{\pm1.0}$ | $82.1^{\pm0.5}$ | $59.8^{\pm1.3}$ | $50.7^{\pm1.1}$ |
| 8-mer | $81.8^{\pm0.9}$ | $77.2^{\pm1.0}$ | $78.1^{\pm0.7}$ | $66.8^{\pm0.3}$ | $60.5^{\pm1.5}$ |
| 10-mer | $75.9^{\pm0.7}$ | $72.5^{\pm1.2}$ | $73.5^{\pm0.8}$ | $67.4^{\pm0.5}$ | $66.2^{\pm1.1}$ |
| 12-mer | $70.8^{\pm1.0}$ | $68.1^{\pm0.6}$ | $68.7^{\pm0.4}$ | $65.9^{\pm0.7}$ | $64.8^{\pm1.5}$ |
| 14-mer | $65.6^{\pm1.6}$ | $62.6^{\pm1.0}$ | $64.2^{\pm1.2}$ | $61.3^{\pm2.4}$ | $61.6^{\pm1.1}$ |
| 16-mer | $61.7^{\pm1.2}$ | $57.9^{\pm0.9}$ | $59.2^{\pm1.0}$ | $61.0^{\pm0.9}$ | $61.1^{\pm1.1}$ |
| poisson-3 | $85.7^{\pm0.7}$ | $68.3^{\pm0.7}$ | $59.2^{\pm0.2}$ | $34.7^{\pm1.8}$ | $26.0^{\pm1.2}$ |
| poisson-5 | $78.6^{\pm0.5}$ | $63.5^{\pm1.0}$ | $60.3^{\pm1.9}$ | $40.3^{\pm1.3}$ | $34.0^{\pm0.7}$ |
| poisson-7 | $68.4^{\pm3.2}$ | $56.6^{\pm0.9}$ | $55.3^{\pm1.0}$ | $40.7^{\pm0.6}$ | $37.0^{\pm0.9}$ |

**D Distributional approach: antiphons**

| | pitch | dep. interval | indep. interval | dep. contour | indep. contour |
|---|---|---|---|---|---|
| neumes | $91.8^{\pm0.5}$ | $79.7^{\pm0.6}$ | $48.4^{\pm0.6}$ | $39.4^{\pm0.8}$ | $30.0^{\pm0.7}$ |
| syllables | $92.1^{\pm0.5}$ | $80.9^{\pm0.7}$ | $52.3^{\pm0.9}$ | $44.4^{\pm0.8}$ | $34.6^{\pm0.8}$ |
| words | $\underline{\mathbf{94.9}^{\pm0.3}}$ | $\underline{\mathbf{92.3}^{\pm0.4}}$ | $\underline{\mathbf{90.1}^{\pm0.6}}$ | $\underline{\mathbf{85.1}^{\pm0.4}}$ | $\underline{\mathbf{82.7}^{\pm0.2}}$ |
| 1-mer | $88.6^{\pm0.2}$ | $54.0^{\pm0.7}$ | $12.4^{\pm0.3}$ | $23.3^{\pm0.6}$ | $12.4^{\pm0.3}$ |
| 2-mer | $92.4^{\pm0.4}$ | $78.2^{\pm0.6}$ | $38.1^{\pm0.9}$ | $29.1^{\pm0.5}$ | $21.4^{\pm0.3}$ |
| 3-mer | $93.7^{\pm0.2}$ | $87.4^{\pm0.3}$ | $69.6^{\pm0.6}$ | $38.0^{\pm0.6}$ | $26.5^{\pm0.3}$ |
| 4-mer | $94.3^{\pm0.4}$ | $90.3^{\pm0.6}$ | $83.5^{\pm0.4}$ | $50.3^{\pm0.8}$ | $33.9^{\pm0.5}$ |
| 5-mer | $93.8^{\pm0.2}$ | $91.2^{\pm0.2}$ | $87.8^{\pm0.3}$ | $62.7^{\pm0.4}$ | $46.0^{\pm0.8}$ |
| 6-mer | $92.7^{\pm0.3}$ | $89.9^{\pm0.3}$ | $89.7^{\pm0.5}$ | $69.0^{\pm0.7}$ | $58.4^{\pm0.5}$ |
| 8-mer | $87.6^{\pm0.4}$ | $85.3^{\pm0.3}$ | $85.3^{\pm0.5}$ | $76.0^{\pm0.8}$ | $70.3^{\pm0.9}$ |
| 10-mer | $81.2^{\pm0.3}$ | $78.1^{\pm0.6}$ | $78.6^{\pm0.8}$ | $72.6^{\pm0.4}$ | $72.0^{\pm0.6}$ |
| 12-mer | $73.8^{\pm0.5}$ | $71.5^{\pm0.7}$ | $72.3^{\pm0.3}$ | $68.8^{\pm0.5}$ | $68.8^{\pm0.8}$ |
| 14-mer | $66.1^{\pm0.4}$ | $63.5^{\pm0.7}$ | $64.2^{\pm0.5}$ | $64.1^{\pm0.8}$ | $63.3^{\pm0.9}$ |
| 16-mer | $60.1^{\pm0.5}$ | $57.0^{\pm0.8}$ | $57.3^{\pm0.5}$ | $57.5^{\pm0.6}$ | $58.5^{\pm0.8}$ |
| poisson-3 | $88.7^{\pm0.4}$ | $77.8^{\pm0.8}$ | $66.0^{\pm1.5}$ | $37.8^{\pm0.6}$ | $29.4^{\pm1.5}$ |
| poisson-5 | $84.6^{\pm0.5}$ | $76.9^{\pm0.9}$ | $72.2^{\pm0.5}$ | $48.3^{\pm1.2}$ | $40.6^{\pm0.6}$ |
| poisson-7 | $79.1^{\pm0.7}$ | $72.4^{\pm0.7}$ | $68.6^{\pm0.8}$ | $52.8^{\pm0.7}$ | $47.6^{\pm1.2}$ |

**Figure S12**. **Classification results with standard deviation.** This is essentially the same figure as figure 5 but now with the mean $F_1$-score $\mu$ and its standard deviation $\sigma$ shown as $\mu^{\pm\sigma}$, computed from five independent runs of the experiment.

**A Classical approach**

| | responsory | antiphon |
|---|---|---|
| final | $38.1^{\pm1.8}$ | $47.8^{\pm0.7}$ |
| ambitus | $48.0^{\pm2.6}$ | $56.9^{\pm1.8}$ |
| initial | $38.2^{\pm1.5}$ | $43.0^{\pm1.3}$ |
| final ambitus | $82.4^{\pm2.0}$ | $76.4^{\pm1.2}$ |
| final initial | $67.1^{\pm2.0}$ | $70.4^{\pm0.7}$ |
| ambitus initial | $62.2^{\pm2.4}$ | $69.7^{\pm2.1}$ |
| final ambitus initial | **$83.3^{\pm1.9}$** | **$82.6^{\pm1.5}$** |

**B Profile approach**

| | responsory | antiphon |
|---|---|---|
| pitch class profile | $76.5^{\pm1.2}$ | $84.5^{\pm0.9}$ |
| pitch profile | **$78.1^{\pm1.1}$** | **$85.5^{\pm1.0}$** |
| repetition profile | $71.3^{\pm1.0}$ | $79.6^{\pm1.5}$ |

**C Distributional approach: responsories**

| | pitch | dep. interval | indep. interval | dep. contour | indep. contour |
|---|---|---|---|---|---|
| neumes | $82.1^{\pm1.0}$ | $66.4^{\pm0.8}$ | $61.9^{\pm0.8}$ | $46.5^{\pm1.9}$ | $41.7^{\pm1.3}$ |
| syllables | $82.4^{\pm1.1}$ | $67.8^{\pm1.0}$ | $64.6^{\pm1.4}$ | $55.6^{\pm1.5}$ | $54.1^{\pm1.3}$ |
| words | $68.5^{\pm2.1}$ | $57.0^{\pm1.5}$ | $55.9^{\pm1.7}$ | $45.5^{\pm2.0}$ | $45.1^{\pm1.4}$ |
| 1-mer | $80.4^{\pm1.8}$ | $44.7^{\pm2.9}$ | $6.8^{\pm0.7}$ | $16.8^{\pm1.8}$ | $6.8^{\pm0.7}$ |
| 2-mer | $82.9^{\pm1.5}$ | $58.8^{\pm2.0}$ | $31.2^{\pm2.2}$ | $22.6^{\pm1.5}$ | $11.1^{\pm2.8}$ |
| 3-mer | $81.3^{\pm1.9}$ | $63.8^{\pm1.8}$ | $53.4^{\pm2.5}$ | $31.1^{\pm2.0}$ | $17.1^{\pm0.6}$ |
| 4-mer | $79.7^{\pm0.8}$ | $62.6^{\pm1.1}$ | $58.4^{\pm1.3}$ | $33.2^{\pm1.1}$ | $27.0^{\pm2.5}$ |
| 5-mer | $75.9^{\pm1.0}$ | $61.5^{\pm1.6}$ | $60.0^{\pm1.9}$ | $35.9^{\pm1.7}$ | $30.0^{\pm1.1}$ |
| 6-mer | $70.3^{\pm1.9}$ | $58.7^{\pm2.0}$ | $59.3^{\pm0.8}$ | $39.0^{\pm2.0}$ | $34.0^{\pm1.2}$ |
| 8-mer | $59.8^{\pm1.6}$ | $51.8^{\pm1.8}$ | $53.3^{\pm2.2}$ | $39.5^{\pm0.8}$ | $36.4^{\pm1.0}$ |
| 10-mer | $51.7^{\pm1.7}$ | $43.5^{\pm1.3}$ | $46.9^{\pm1.7}$ | $39.5^{\pm1.1}$ | $38.2^{\pm2.9}$ |
| 12-mer | $45.2^{\pm1.9}$ | $40.1^{\pm1.4}$ | $41.2^{\pm2.6}$ | $37.6^{\pm1.9}$ | $36.3^{\pm2.9}$ |
| 14-mer | $39.5^{\pm0.7}$ | $34.9^{\pm2.4}$ | $36.9^{\pm1.6}$ | $33.1^{\pm2.0}$ | $34.3^{\pm1.4}$ |
| 16-mer | $36.8^{\pm2.0}$ | $29.9^{\pm1.2}$ | $31.5^{\pm1.7}$ | $30.3^{\pm1.3}$ | $32.4^{\pm1.1}$ |
| poisson-3 | $73.7^{\pm2.3}$ | $49.2^{\pm1.1}$ | $42.6^{\pm1.8}$ | $23.1^{\pm3.0}$ | $16.4^{\pm2.8}$ |
| poisson-5 | $62.5^{\pm1.1}$ | $43.3^{\pm1.6}$ | $41.8^{\pm0.8}$ | $23.6^{\pm2.5}$ | $20.6^{\pm4.7}$ |
| poisson-7 | $53.0^{\pm2.2}$ | $38.3^{\pm3.0}$ | $37.1^{\pm1.6}$ | $25.2^{\pm2.1}$ | $18.8^{\pm3.3}$ |

**D Distributional approach: antiphons**

| | pitch | dep. interval | indep. interval | dep. contour | indep. contour |
|---|---|---|---|---|---|
| neumes | $88.0^{\pm0.9}$ | $71.6^{\pm0.8}$ | $42.0^{\pm0.9}$ | $34.4^{\pm1.6}$ | $25.4^{\pm0.8}$ |
| syllables | $88.0^{\pm0.7}$ | $72.4^{\pm0.6}$ | $44.4^{\pm1.1}$ | $38.0^{\pm0.5}$ | $29.0^{\pm1.4}$ |
| words | $88.7^{\pm0.4}$ | $83.2^{\pm1.2}$ | $82.5^{\pm0.8}$ | $77.2^{\pm1.1}$ | $75.5^{\pm1.1}$ |
| 1-mer | $85.9^{\pm1.6}$ | $49.7^{\pm1.1}$ | $9.0^{\pm0.8}$ | $19.8^{\pm0.9}$ | $9.4^{\pm0.5}$ |
| 2-mer | $88.9^{\pm1.2}$ | $73.2^{\pm0.7}$ | $34.4^{\pm1.3}$ | $23.9^{\pm1.7}$ | $17.4^{\pm0.6}$ |
| 3-mer | $89.2^{\pm1.0}$ | $80.9^{\pm0.7}$ | $63.8^{\pm1.0}$ | $34.0^{\pm1.9}$ | $21.6^{\pm0.6}$ |
| 4-mer | $89.9^{\pm0.6}$ | $83.5^{\pm1.1}$ | $76.7^{\pm1.3}$ | $43.7^{\pm1.5}$ | $29.4^{\pm1.1}$ |
| 5-mer | $88.5^{\pm0.9}$ | $82.8^{\pm0.5}$ | $80.1^{\pm0.5}$ | $52.4^{\pm2.1}$ | $40.4^{\pm1.2}$ |
| 6-mer | $86.4^{\pm0.5}$ | $80.9^{\pm1.1}$ | $79.9^{\pm0.7}$ | $57.0^{\pm1.6}$ | $48.1^{\pm1.1}$ |
| 8-mer | $78.4^{\pm1.1}$ | $73.8^{\pm1.7}$ | $74.9^{\pm1.6}$ | $59.4^{\pm1.7}$ | $55.2^{\pm1.1}$ |
| 10-mer | $69.0^{\pm2.3}$ | $62.3^{\pm1.3}$ | $63.1^{\pm0.8}$ | $54.9^{\pm2.3}$ | $54.1^{\pm1.6}$ |
| 12-mer | $58.8^{\pm0.5}$ | $53.3^{\pm0.6}$ | $53.4^{\pm0.7}$ | $50.2^{\pm0.8}$ | $49.3^{\pm1.0}$ |
| 14-mer | $50.4^{\pm0.7}$ | $44.8^{\pm1.1}$ | $45.2^{\pm1.4}$ | $43.7^{\pm2.2}$ | $43.7^{\pm2.0}$ |
| 16-mer | $47.5^{\pm1.3}$ | $42.2^{\pm1.6}$ | $41.2^{\pm1.8}$ | $38.2^{\pm1.1}$ | $39.2^{\pm2.0}$ |
| poisson-3 | $82.5^{\pm1.3}$ | $67.9^{\pm1.1}$ | $55.6^{\pm1.5}$ | $28.2^{\pm1.2}$ | $20.2^{\pm1.1}$ |
| poisson-5 | $75.8^{\pm1.1}$ | $64.0^{\pm1.0}$ | $57.7^{\pm1.2}$ | $34.1^{\pm1.8}$ | $29.3^{\pm1.8}$ |
| poisson-7 | $68.7^{\pm1.4}$ | $59.6^{\pm1.4}$ | $56.5^{\pm0.6}$ | $39.2^{\pm2.1}$ | $34.2^{\pm1.6}$ |

**Figure S13**. **Main results on subset.** Cantus often contains several variants of the same melody, as shown in Figure S11. As discussed in the main text, this is a difficult issue that for example also applies to the Essen folk-song collection. We decided to repeat our experiments on a subset of the data where we excluded melody variants. We heuristically identified melody variants by randomly picking one chant from all sets of chants that have the same Cantus ID and mode. This resulted in a set of 1763 responsories and 4159 antiphons. This figure shows the main classification results on this subset of the data. Clearly, the performance of all models drops. The drop is greatest for responsories across models. The main result that only natural units maintain high performance, even on contour representations, nevertheless stand. It should be noted that in antiphons, some *n*-grams now outperform the natural units (by less then 2%) when using pitch and dependent interval representation. This does not reflect a large change in performance: in the main results, these *n*-grams are also deviated from top performance by no more than 2%. These results are further discussed in the main text.