

Learning to classify organic and conventional wheat - a machine-learning driven approach using the MeltDB 2.0 metabolomics analysis platform

Nikolas Kessler, Anja Bonte, Stefan P Albaum, Paul Mäder, Monika Messmer, Alexander Goesmann, Karsten Niehaus, Georg Langenkämper and Tim W Nattkemper

Journal Name:	Frontiers in Bioengineering and Biotechnology
ISSN:	2296-4185
Article type:	Original Research Article
Received on:	04 Dec 2014
Accepted on:	03 Mar 2015
Provisional PDF published on:	03 Mar 2015
Frontiers website link:	www.frontiersin.org
Citation:	Kessler N, Bonte A, Albaum SP, Mäder P, Messmer M, Goesmann A, Niehaus K, Langenkämper G and Nattkemper TW(2015) Learning to classify organic and conventional wheat - a machine-learning driven approach using the MeltDB 2.0 metabolomics analysis platform. <i>Front. Bioeng. Biotechnol.</i> 3:35. doi:10.3389/fbioe.2015.00035
Copyright statement:	© 2015 Kessler, Bonte, Albaum, Mäder, Messmer, Goesmann, Niehaus, Langenkämper and Nattkemper. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY) . The use, distribution and reproduction in other forums is permitted, provided the original author(s) or licensor are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

This Provisional PDF corresponds to the article as it appeared upon acceptance, after rigorous peer-review. Fully formatted PDF and full text (HTML) versions will be made available soon.

Learning to classify organic and conventional wheat - a machine-learning driven approach using the MeltDB 2.0 metabolomics analysis platform

Nikolas Kessler^{1,2}, Anja Bonte³, Stefan P. Albaum², Paul Mäder⁴, Monika Messmer⁵, Alexander Goesmann⁶, Karsten Niehaus⁷, Georg Langenkämper³, and Tim W. Nattkemper^{1*}

¹*Biodata Mining Group, Faculty of Technology, Bielefeld University, Bielefeld, Germany*

²*Bioinformatics Resource Facility, Center for Biotechnology, Bielefeld University, Bielefeld, Germany*

³*Department of Safety and Quality of Cereals, Max Rubner-Institut, Detmold, Germany*

⁴*Department of Soil Sciences, Research Institute of Organic Agriculture (FiBL), Frick, Switzerland*

⁵*Department of Crop Sciences, Research Institute of Organic Agriculture (FiBL), Frick, Switzerland*

⁶*Bioinformatics and Systems Biology, Justus-Liebig-University Gießen, Gießen, Germany*

⁷*Department of Proteome and Metabolome Research, Center for Biotechnology, Bielefeld University, Bielefeld, Germany*

Correspondence*:

Tim W. Nattkemper

Biodata Mining and Applied Neuroinformatics Group, Faculty of Technology, Bielefeld University, Universitätsstr. 25, 33615 Bielefeld, Germany, tim.nattkemper@uni-bielefeld.de

2 ABSTRACT

3 We present results of our machine learning approach to the problem of classifying GC-MS
4 data originating from wheat grains of different farming systems. The aim is to investigate the
5 potential of learning algorithms to classify GC-MS data to be either from conventionally grown or
6 from organically grown samples and considering different cultivars. The motivation of our work
7 is rather obvious nowadays: increased demand for organic food in post-industrialized societies
8 and the necessity to prove organic food authenticity. The background of our data set is given
9 by up to eleven wheat cultivars that have been cultivated in both farming systems, organic and
10 conventional, throughout three years. More than 300 GC-MS measurements were recorded and
11 subsequently processed and analyzed in the MeltDB 2.0 metabolomics analysis platform, being
12 briefly outlined in this paper. We further describe how unsupervised (t-SNE, PCA) and supe-
13 rvised (SVM) methods can be applied for sample visualization and classification. Our results
14 clearly show that years have most and wheat cultivars have second-most influence on the meta-
15 bolic composition of a sample. We can also show, that for a given year and cultivar, organic and
16 conventional cultivation can be distinguished by machine-learning algorithms.

17

18 **Keywords:** Metabolome informatics, statistics, metabolomics, computational metabolomics, organic farming, food authentication,
19 machine learning

1 INTRODUCTION

20 The increasing awareness of the benefits of healthy eating has tremendously risen the popularity of organic
21 food - a development that was not least stirred up by the manifold food scandals grabbing the headlines
22 in recent years. Directly resulting from this popularity but in particular from organic food's great market
23 potential, there emerged a significant interest in the authenticity of food declared as organic [Capuano
24 et al., 2013]. Metabolomics technologies have proven successful for several task of food authentication
25 [Cubero-Leon et al., 2014]. In this study, we investigate the potential of metabolomics profiling techni-
26 ques, bioinformatics and machine learning to distinguish organically grown wheat from conventionally
27 grown wheat. To this end, a total of more than 300 gas chromatography-mass spectrometry (GC-MS)
28 measurements from both types of treatments were recorded and analyzed. Samples comprised eleven dif-
29 ferent wheat cultivars from up to three different years, obtained from the DOK field trial in Switzerland
30 [Mäder et al., 2002]. This comprehensive field trial compared organic and conventional farming systems,
31 using strictly controlled conditions. In previous work [Bonte et al., 2014], we already presented meta-
32 bolite profiling data obtained from the DOK wheat samples of the harvest year 2007. Röhlig and Engel
33 [2010] have applied principal component analysis (PCA) and analysis of variance (ANOVA) to a very
34 similar dataset. In the scope of this work, we substantially extended the DOK data basis from 2007 by
35 additionally analyzing samples from the 2009 and 2010 harvest years. The particular focus of this work
36 was placed on the potential of machine learning methods as tools for automated data classification. Furth-
37 ermore the new approach is metabolite-agnostic: It does not rely on correct metabolite identification and it
38 does not rely on single biomarkers with significant level differences. The latter is a core advantage of this
39 approach, as literature reveals that only slight (not significant) metabolite level changes can be accounted
40 on the farming systems [Röhlig and Engel, 2010; Laursen et al., 2011; Bonte et al., 2014].

41 All GC-MS measurements were automatically preprocessed and then carefully annotated in our MeltDB
42 2.0 metabolomics analysis platform [Neuweger et al., 2008; Kessler et al., 2013]. MeltDB allowed us to
43 apply a well-established routine in high-dimensional molecular data analysis. After pre-processig (peak
44 picking, normalization, profiling etc) the data is represented as a table of dimension $n \times D$, with n =
45 number of samples and D = signal dimension (i.e. the metabolic profile). The first aim is to search for
46 hidden regularities, relationships, and correlations in the data. To this end, unsupervised learning, i.e.
47 dimensional reduction can be applied. Concretely, the two unsupervised methods principal component
48 analysis (PCA) and t-distributed stochastic neighbor embedding (t-SNE, van der Maaten and Hinton
49 [2008]) were used to investigate the inter- and intra-class variances in the entire dataset as well as in
50 particular subsets of the data.

51 Second, the data was analyzed towards the question, if it can be classified into distinct semantic catego-
52 ries (like conventional / organic treatment in this case). We therefore applied the two supervised machine
53 learning methods random forests (RF, Breiman [2001]) as well as support vector machines (SVM, Vapnik
54 [1999]). The overall aim was to establish a classifier to distinguish between organic and conventional
55 wheat, despite the influences of the years of growth and different cultivars.

56 In the following, the analytical approach is described in detail, as well as the separation results for the
57 investigated factors treatment, cultivar and year. All presented computational methods were implemented
58 within the MeltDB 2.0 platform and can be applied on other datasets as well.

2 MATERIAL & METHODS

2.1 PLANT MATERIAL

59 Wheat grains of up to 11 different cultivars originated from the DOK (D: bio-dynamic, O: bio-organic,
60 K: 'konventionell' German for conventional, i.e. integrated, farming system) field trial, which is located
61 at Therwil (7°33' E, 47°30' N) close to Basel (Switzerland). Detailed Information on the DOK long-term
62 field trial is given by **Mäder et al.** [2002]. Wheat grains of the cultivar Runal were analyzed from the
63 three harvest years 2007, 2009 and 2010. In 2008 wheat was not grown in the trial. Further, the 10 wheat
64 cultivars 'Rouge de Bordeaux', 'Mont Calme 245', 'Probus', 'CCP' (composite cross-population; for ease
65 of reading CCP is referred to as a cultivar), 'Scaro', 'Sandomir', 'DJ 9714', 'Antonius', 'Caphorn' and
66 'Titlis' were integrated into the wheat plots of the long term trial of the harvest year 2007. In the 2010
67 cultivation period, cultivars 'Mont Calme 245' and 'DJ9714' were not available, leaving the remaining 8
68 cultivars mentioned previously for analysis in this work. A detailed description of the layout and design
69 of the experiment comprising all winter wheat cultivars was published [**Hildermann et al.**, 2009].

70 Thus only some essential information about the DOK field trial is considered here. The trial comprises
71 several organic and conventional farming systems, each system being repeated in four field plots. The
72 experimental design was a split plot with systems as the main factor and wheat cultivars as the secondary
73 factor. For this work we choose to analyze the two farming systems biodynamic 2 (D), (henceforth, orga-
74 nic) and conventional (M). These two farming systems were quite different with respect to fertilization
75 and further plant treatment (see below), but at the same time were still within the range of standard organic
76 and conventional farming.

77 The organic system received composted manure and slurry at a fertilization level of 1.4 livestock units
78 per hectare, equivalent to 66 kg N(total) ha⁻¹. Fertilization in the conventional system was done exclusi-
79 vely with mineral fertilizer at 140 kg N(total) ha⁻¹. Both farming systems also differed in plant protection
80 practice. The conventional system followed the guidelines of integrated farming, using fungicides, inse-
81 cticides and herbicides only if needed. The biodynamic farming allowed only mechanical plant treatments
82 and indirect methods to control weeds, pests and diseases. Grains of both farming systems were harve-
83 sted when completely ripe, with moisture content below 140 g kg⁻¹. Of each of the four individual field
84 plots per agricultural system, one sample was taken for each cultivar and farming system. Before further
85 experimental usage, grain material was stored at a constant temperature of 18 °C.

2.2 WHEAT SAMPLE PREPARATION AND GC-MS ANALYSIS

86 Cleaning of wheat samples from impurities and broken grains, grain storage, grinding and extraction as
87 well as measurement of metabolites using GC-MS analysis was exactly performed as described by **Bonte**
88 **et al.** [2014].

2.3 DATA PROCESSING IN MELTDB 2.0

89 All data gained by GC-MS analysis were preprocessed and annotated within the MeltDB 2.0 metabolo-
90 mics software platform. Peaks were obtained using the Warped Peak Detection tool. Retention indices
91 were obtained semi-automatically, using MeltDB's RISimple tool and a manually defined list of expe-
92 cted retention times for each batch of measurements. Next, a profiling was run to annotate peaks that
93 are common throughout multiple chromatograms, i.e. they have a similar retention index and a similar
94 EIC spectrum. Similarly, all chromatograms were matched against reference spectra to annotate peaks
95 as identified compounds where possible. However, the subsequent approach does not rely on the iden-
96 tification of compounds, but rather uses it to limit the feature space to molecules of potential biological
97 interest. The parameterizations for these processing tools can be found in Table 1.

98 Results from automated metabolite identification were revised manually to discard erroneous annotati-
99 ons, but also to manually create annotations that were missed in a few chromatograms only. Peaks that

Table 1 - Supporting Information. Parameters that were applied for preprocessing tools.

Tool	Description	Parameter	Value
Warped Peak Detection	Mexican-wavelet based peak detection, which can be rerun locally (at certain RT).	FWHM SN	7 10
RISimple	Detects and tags retention indices based on heir characteristic spectra.	ion filter	57, 71, 85, 99
Multiple Profiling	Gives peaks across chromatograms a common TAG if they are similar.	Retention time window	20-35 s
Reference List	Annotates peaks that match reference spectra, uses dot-product.	RT Window	20 s

Table 2. Number of samples for each combination of factors 'farming system', 'year', and 'cultivar'.

Farming system	Year	Cultivar											Σ
		Antonius	Caphorn	CCP	DJ 9714	Mont Calme 245	Probus	Rouge de Bordeaux	Runal	Sandomir	Scaro	Titlis	
conventional organic	2007	7 8	8 7	7 7	7 7	8 7	7 8	7 8	8 7	6 8	6 7	8 7	160
conventional organic	2009								8 8				16
conventional organic	2010	8 8	8 8	7 8			8 8	8 8	7 7	8 7	7 7	8 7	137

100 were missed in a minority of samples only, were requantified using the Warped Peak Detection tool.
101 Subsequently all data in the obtained feature table was centered and scaled using R.

102 In total 313 samples were analyzed. From the years 2007, 2009, and 2010 these comprise 160, 16, and
103 137 samples, respectively. How many cultivars were available in each year is mentioned in section 2.1.
104 For each combination of year and cultivar, 13 to 16 samples comprising duplicate metabolite extractions
105 for grains from most field plots, were analyzed. From these one half was treated organically, and the other
106 half was treated conventionally. A detailed listing of all samples is given in Tab. 2.

2.4 UNSUPERVISED LEARNING / DIMENSIONAL REDUCTION

107 The result of the pre-processing step is an $n \times D$ data table ($n=313$, total number of samples (years,
108 cultivars and treatments combined); $D = 36$, number of compounds consistently annotated in all sam-
109 ples). Additionally, subsets of this data table, for example all samples from only one year or cultivar, have
110 been analyzed too. Although it may be tempting to instantly apply supervised learning to the problem of

111 classifying the data rows into conventional and biological treatment, we first applied some information
112 visualization in advance to avoid unpleasant black box effects and to gain a mental model of the data.
113 Information Visualisation uses different data displays which are inspected by human experts to under-
114 stand the data or to build hypotheses for the hidden structures in the data. These are the foundation for any
115 subsequent attempt to apply supervised learning. We propose to inspect displays obtained with two diffe-
116 rent dimensional reduction techniques. First, we applied PCA since this is a well established statistics tool
117 in high dimensional data analytics and is fully sufficient to understand data with a linear sub-structure.
118 Since data stemming from systems biology experiments can not be expected to have such intrinsic linear
119 structure we used another method which has been proposed in the field of machine learning, the t-SNE.
120 In several real world applications for computational biology [Abdelmoula et al., 2014; Bushati et al.,
121 2011; Jamieson et al., 2010] t-SNE has shown to be capable of projecting non-linear data structure while
122 well preserving the local features (i.e. neighborhoods) of the data.

123 The dimensional reductions were performed using the R statistical software [R_Development_Core_Team,
124 2011] and the “tsne” package by Donaldson [2012].

2.5 SUPERVISED LEARNING / CLASSIFICATION

125 The same $n \times D$ data table was used to explore whether a machine learning algorithm such as the Support
126 Vector Machine (SVM, Vapnik [1999]) with a polynomial kernel [Karatzoglou et al., 2004] can be
127 trained to classify the data rows into conventional and biological treatment. In a first step for each subset
128 (e.g. data from one year only), the machine learning algorithm was trained and tested on 80 percent
129 of the data (randomly selected). Afterwards, the remaining 20 percent out-of-the-bag data was used for
130 validation, i.e. to finally evaluate the performance of the classifier constructed using the 80 percent of the
131 data. To train and optimize the SVM, a parameter tuning was performed using a 25-fold resampling for
132 Leave-Group-Out-Cross-Validation (LGOCV) of the training partition. For this cross validation, again 75
133 percent of the training partition were used for training, and 25 percent were used for validation in each
134 iteration. The best set of parameters, that led to the best accuracy according to the LGOCV, was then once
135 more validated on the 20 percent of the data that was kept back initially. The classification results on these
136 latter 20 percent were evaluated in a confusion matrix to infer the accuracy of the trained SVM.

137 Using the very same subsets and partitions random forest (RF, Liaw and Wiener [2004]) was performed
138 as well. The RF training was done with a 20-fold resampling and a parameter tune length of twelve.

139 Supervised machine learning methods were performed in R as well, using the “caret” package [Kuhn,
140 2008; Kuhn et al., 2008].

3 RESULTS

3.1 UNSUPERVISED LEARNING / DIMENSIONAL REDUCTION

141 Both, PCA (in the first two principal components, see Fig. 1) and t-SNE, are capable of separating the
142 presented wheat samples into clusters according to the factor year. Within one year the PCA will group
143 samples according to cultivars, though with considerable overlap as can be seen in Fig. 2. Conversely,
144 within one cultivar samples will be grouped according to the year (see Fig. 3). When one year and one
145 cultivar are investigated in any combination, all data typically clusters into the two groups representing
146 either dynamically or conventionally grown wheat. Most of these clusters show at least some overlap
147 though. The t-SNE method is less applicable to smaller datasets and thus was applied to the complete data
148 table only. Fig. 5 shows how t-SNE groups all samples by year at first, and then into subclusters according
149 to their cultivars. The latter subclusters themselves are again split in two groups each, which correspond to
150 the two farming systems, as can be seen in Fig. 6. Fig. 5 and 6 again visualise strikingly how the metabolic
151 profile is mainly influenced by year, then by cultivar, and at least by the farming systems.

152 Nevertheless it is still obvious that the treatment caused measurable differences in the metabolic compo-
153 sition of the wheat samples. The first two principal components of the PCA in Fig. 3 already reveal, that
154 the clusters for the years 2007, 2009 and 2010 are split in themselves to form subclusters of conventio-
155 nally and organically grown wheat. This points out that later principal components with different loadings
156 may expose structures in the data that are mainly based on the factor treatment. Fig. 4 plots the second
157 and fourth principal components of the PCA that has already been introduced in Fig. 3. In this case it is
158 clearly visualised how the fourth principal component can be used to separate the samples according to
159 the levels organic and conventional.

3.2 SUPERVISED LEARNING / CLASSIFICATION

160 The results from the PCAs revealed that there are structures in the data that allow for a separation of
161 conventionally and organically grown wheat. Even though the main clustering is driven by factor year,
162 these clusters still form subclusters according to cultivar, which again are clustered by the two farming
163 systems. These substructures suggest, that SVMs can be constructed to win classifiers for the problem.
164 In fact, SVMs trained and tested on the entire dataset (all years, all cultivars, both treatments) to classify
165 by treatment reached an accuracy of 0.9032 (p-value = $1.486e-11$, see Tab. 3) on the validation set. Even
166 better accuracies can be observed when investigating subsets of the data (for example accuracy = 0.9677,
167 p-value = $3.746e-08$ within year 2007). But the smaller the subsets, the smaller the testing partitions, the
168 less representative are any outcomes. Thus we will not trust the classifiers for in-cultivar or even in-year-
169 and-cultivar problems to be flawless, even though in these cases accuracy values may approximate one
170 easily.

171 The interesting question would be, if it is possible to obtain such a trained classifier from a number of
172 (past) years, that can then be applied to classify samples from another (e.g. the present) year. This however
173 turns out to be not possible on the basis of the available data from the three growing seasons. For example,
174 when a SVM, trained on data from 2007, is applied to classify data from 2010 it performs with an accuracy
175 of 0.5547, which is hardly favourable to plain guesses. The reason for this poor performance seems to be
176 the massive influence of the seasonal conditions, i.e. the factor 'year'. This calls for continuing research
177 using more samples from more years and cultivars to cover the molecular variance more appropriately.
178 Estimations on the variable importance [Kuhn et al., 2008] for the three years were calculated based on
179 the SVM results and added to the supplemental information. Here it is striking that e.g. *myo*-inositol,
180 which has previously been reported as a potential marker for farming systems [Röhlig and Engel, 2010;
181 Bonte et al., 2014], was most important for classification in 2007 but almost least important in 2009 and
182 2010. Such inhomogeneous variable importances additionally suggest a year-by-year strategy for training
183 and classification.

184 Table 3 summarizes the SVM results. Please note, that classification results for year 2009 are not repor-
185 ted here: with only one cultivar (Runal) and thus only 16 samples the subset is too small to generate
186 reliable results. The 2009 samples are part of the analysis of the entire dataset, though.

187 Overall random forest (RF) as described in the methods section led to similar classification results,
188 but showed slightly lower accuracies in the in-year subsets. Thus no detailed results are shown in the
189 manuscript. However, we explicitly do not suggest to ignore random forest as a potential alternative for
190 support vector machines in this scope.

4 DISCUSSION

191 The main goal of this study was to investigate whether a classification of organically and conventionally
192 grown wheat can be done, based on GC-MS metabolite measurements of wheat grains from different
193 years and cultivars. Results from the unsupervised machine learning methods PCA and t-SNE show, that
194 the strongest variation in the data can be found in samples from different years. This may be in part
195 due to different environmental influences and also due to systematic errors that inevitably will occur in

Table 3. Results of the Support Vector Machines, trained and tested on different subsets of all samples. Measures are given for the evaluation results and are based on the confusion matrix for classification as biological or conventional farming system.

Trained on	Tested on	n_{Test}	Accuracy	NIR ¹	p-value ²	Sensitivity	Specificity	PPV ³	NPV ⁴
2007	2007	31	0.9677	0.52	3.75E-08	1	0.9375	0.9375	1
2010	2010	26	0.8846	0.5	4.40E-05	0.9231	0.8462	0.8571	0.9167
2007	2010	137	0.5547	0.5	0.1333	0.2754	0.8382	0.6333	0.5327
2010	2007	160	0.5562	0.51	0.1177	0.8101	0.3086	0.5333	0.625
2007, 09, 10	2007, 09, 10	62	0.9032	0.5	1.49E-11	0.9032	0.9032	0.9032	0.9032

¹No information rate: the larger class percentage; ²Exact binomial test [Accuracy > NIR];

³Positive predictive value; ⁴Negative predictive value

196 analyses from different years. Other studies report on the same obstructive effects [Röhlig and Engel,
197 2010; Laursen et al., 2011]. On the other hand though, this allows to extend the data basis every year.
198 This demands robust classifiers that are able to cope with these kinds of problems, besides “distracting”
199 factors like year and cultivar. Further studies will additionally have to consider geographical influences on
200 the metabolic composition of wheat grains.

201 Peaks from all 313 samples have been carefully annotated to achieve 36 consistently quantified featu-
202 res throughout the entire data set. These have first been explored with dimensional reduction methods
203 like PCA and t-SNE to find the predominant structures in the data table. Then, supervised machine
204 learning methods have been trained and applied to investigate in how far classifiers for organically and
205 conventionally grown wheat can be created.

206 The considerably strong differences in samples from different years make it impossible though, to apply
207 a classifier that was trained using data from year a_1 to distinguish data from another year a_2 . To create a
208 classifier for any year a_x , data from this a_x must be part of the training data set. PC analyses also suggest
209 that it will be beneficial to concentrate on one cultivar or to have a broad data basis of many cultivars to
210 cover variances that derive from this factor.

211 Support vector machines trained and applied on all samples from the same year, as well as SVMs trained
212 and tested on all years, performed with high accuracies above or close to 0.9. This clearly outperforms the
213 ability of PCA to separate samples according to the applied farming system, unless samples derive from
214 the same cultivar. For comparison, we also performed a study using Random Forests [Breiman, 2001]
215 instead of SVMs for classification. Random Forests (RF) have the advantage to be much faster and more
216 efficient than SVMs and they have the potential to offer some insight into the semantics of the decision
217 function, but the parameters are more difficult to optimize. However, the classification performances were
218 only slightly different from those obtained with SVMs and inferior in in-year analyses, so we did not
219 include those in the manuscript.

220 The here presented machine learning tools are not meant to substitute traditional statistical methods,
221 such as ANOVA, but provide a metabolite-agnostic approach for sample classification where reliable bio-
222 markers are not known. Additionally, they may contribute a starting point for focused statistical analyses
223 of single compounds that appear promising according to the computed variable importance estimations.

224 An analytical approach that aims more for specific compounds as biological markers can be found in
225 the publication of Bonte et al. [2014], where more traditional statistical methods have been applied.
226 The methods presented in the manuscript at hand do not depend on the identification of compounds or

227 the determination of the biological meaning of any features. The approach rather relies on a consistently annotated data set. Nevertheless it is constructive to do compound identification to be able to base further biomarker research on these studies. Additionally, reducing the feature set to verified biological compounds minimizes the risk of systematic errors through background noise. Variable importance estimations based on the SVM results of the three years have thus been added to the supplemental information. The integration of the discussed approaches might finally lead to a set of metabolites that can be used as reliable biomarkers for conventional or biodynamic farming systems.

DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

234 The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

236 Design and implementation of the DOK experiment: PM, CT. Conception and design of wet-lab experiments: AB, KN, GL. Conception and design of data analysis: NK, SPA, AG, KN, TWN. GC-MS measurements: AB. Data processing: NK, AB. Data analysis and software implementation: NK. Contribution of raw data/materials/analysis tools: SPA, PM, CT, AG, KN, GL, TWN. Preparation of the manuscript: NK, AB, SPA, PM, CT, AG, KN, GL, TWN.

ACKNOWLEDGEMENT

241 We thank Agroscope, Institute for Sustainability Sciences (ISS) (Switzerland) for assistance in field techniques and we thank the BRF support team for expert IT support.

243 *Funding:* The authors acknowledge the financial support provided by the German Federal Ministry of Food and Agriculture under the Federal Scheme for Organic Farming and other forms of sustainable agriculture (Project 08OE023). Wheat cultivar testing in the DOK trial was supported within the FP 7 project NUE-CROPS (EU-FP7 222-645). NK was supported by a fellowship of the CLIB Graduate Cluster Industrial Biotechnology.

REFERENCES

- 248 Abdelmoula, W. M., Škrášková, K., Balluff, B., Carreira, R. J., Tolner, E. a., Lelieveldt, B. P. F., et al.
249 (2014), Automatic generic registration of mass spectrometry imaging data to histology using nonlinear
250 stochastic embedding., *Analytical chemistry*, 86, 18, 9204–11, doi:10.1021/ac502170f
- 251 Bonte, A., Neuweger, H., Goesmann, A., Thonar, C., Mäder, P., Langenkämper, G., et al. (2014), Metabo-
252 lite profiling on wheat grain to enable a distinction of samples from organic and conventional farming
253 systems., *Journal of the science of food and agriculture*, , January, doi:10.1002/jsfa.6566
- 254 Breiman, L. (2001), Random forests, *Machine learning*, 45, 1, 5–32
- 255 Bushati, N., Smith, J., Briscoe, J., and Watkins, C. (2011), An intuitive graphical visualization technique
256 for the interrogation of transcriptome data., *Nucleic acids research*, 39, 17, 7380–9, doi:10.1093/nar/
257 gkr462
- 258 Capuano, E., Boerrigter-Eenling, R., van der Veer, G., and van Ruth, S. M. (2013), Analytical authentication
259 of organic products: an overview of markers, *Journal of the Science of Food and Agriculture*, 93,
260 October, 12–28, doi:10.1002/jsfa.5914

- 261 Cubero-Leon, E., Peñalver, R., and Maquet, A. (2014), Review on metabolomics for food authentication,
262 *Food Research International*, 60, 95–107, doi:10.1016/j.foodres.2013.11.041
- 263 Donaldson, J. (2012), tsne: T-distributed Stochastic Neighbor Embedding for R (t-SNE)
- 264 Hildermann, I., Thommen, A., Dubois, D., Boller, T., Wiemken, A., and Mäder, P. (2009), Yield and
265 baking quality of winter wheat cultivars in different farming systems of the DOK long-term trial,
266 *Journal of the Science of Food and Agriculture*, 89, 14, 2477–2491, doi:10.1002/jsfa.3750
- 267 Jamieson, A. R., Giger, M. L., Drukker, K., Li, H., Yuan, Y., and Bhooshan, N. (2010), Exploring
268 nonlinear feature space dimension reduction and data representation in breast CADx with Laplacian
269 eigenmaps and t-SNE, *Medical Physics*, 37, 1, 339, doi:10.1118/1.3267037
- 270 Karatzoglou, A., Smola, A., Hornik, K., and Zeileis, A. (2004), kernlab An S4 Package for Kernel
271 Methods in R, *Journal of Statistical Software*, 11, 9
- 272 Kessler, N., Neuweger, H., Bonte, A., Langenkämper, G., Niehaus, K., Nattkemper, T. W., et al. (2013),
273 MeltDB 2.0-advances of the metabolomics software system., *Bioinformatics (Oxford, England)*, 29,
274 19, 2452–2459, doi:10.1093/bioinformatics/btt414
- 275 Kuhn, M. (2008), Building predictive models in R using the caret package, *Journal of Statistical Software*,
276 28, 5, 1–26
- 277 Kuhn, S., Egert, B., Neumann, S., and Steinbeck, C. (2008), Building blocks for automated elucidation
278 of metabolites: machine learning methods for NMR prediction., *BMC bioinformatics*, 9, 400, doi:10.
279 1186/1471-2105-9-400
- 280 Liaw, A. and Wiener, M. (2002), Classification and Regression by randomForest, *R News*, 2, 3, 18–22
- 281 Laursen, K. H., Schjoerring, J. K., Olesen, J. r. E., Askegaard, M., Halekoh, U., and Husted, S. r. (2011),
282 Multielemental fingerprinting as a tool for authentication of organic wheat, barley, faba bean, and
283 potato, *Journal of Agricultural and Food Chemistry*, 59, 4385–4396, doi:10.1021/jf104928r
- 284 Mäder, P., Fliessbach, A., Dubois, D., Gunst, L., Fried, P., and Niggli, U. (2002), Soil fertility and
285 biodiversity in organic farming., *Science (New York, N.Y.)*, 296, 5573, 1694–7, doi:10.1126/science.
286 1071148
- 287 Neuweger, H., Albaum, S. P., Dondrup, M., Persicke, M., Watt, T., Niehaus, K., et al. (2008), MeltDB:
288 a software platform for the analysis and integration of metabolomics experiment data., *Bioinformatics*
289 (*Oxford, England*), 24, 23, 2726–32, doi:10.1093/bioinformatics/btn452
- 290 R_Development_Core_Team (2011), R: A Language and Environment for Statistical Computing, doi:
291 ISBN3-900051-07-0
- 292 Röhlig, R. M. and Engel, K. H. (2010), Influence of the input system (Conventional versus organic far-
293 ming) on metabolite profiles of maize (*Zea mays*) kernels, *Journal of Agricultural and Food Chemistry*,
294 58, 3022–3030, doi:10.1021/jf904101g
- 295 van der Maaten, L. and Hinton, G. (2008), Visualizing data using t-SNE, *Journal of Machine Learning*
296 *Research*, 9, 85, 2579–2605
- 297 Vapnik, V. N. (1999), An overview of statistical learning theory., *IEEE transactions on neural networks /*
298 *a publication of the IEEE Neural Networks Council*, 10, 5, 988–99, doi:10.1109/72.788640

FIGURES

299 **Figure 1.** The principal component analysis on the entire dataset of all samples throughout all years,
300 cultivars and treatments shows, that the first two components mainly separate samples by the factor year.
301 A separation by the factor farming system is not possible.

302 **Figure 2.** A principal component analysis performed on a dataset from one year only will mainly cluster
303 samples by their cultivar, regardless of the applied farming system. This PCA is based on samples from
304 the year 2007.

305 **Figure 3.** Similar to Fig. 1, in the principal component analysis on a dataset of only one cultivar - here
306 “Runal” is shown - the first principal components separate samples by factor year.

307 **Figure 4.** Plotting samples from one cultivar (here “Runal”) along the principal components two and
308 four shows, that a separation by farming system might be possible even though the main variance is caused
309 by the factor year.

310 **Figure 5.** The t-SNE method applied to all samples results in clusters and sub clusters formed according
311 to the factor year and cultivar, respectively.

312 **Figure 6.** The same t-SNE result as in Fig. 5, but colored by farming system: Clusters representing
313 cultivars form subclusters according to the factor farming system.

Figure 1.TIF

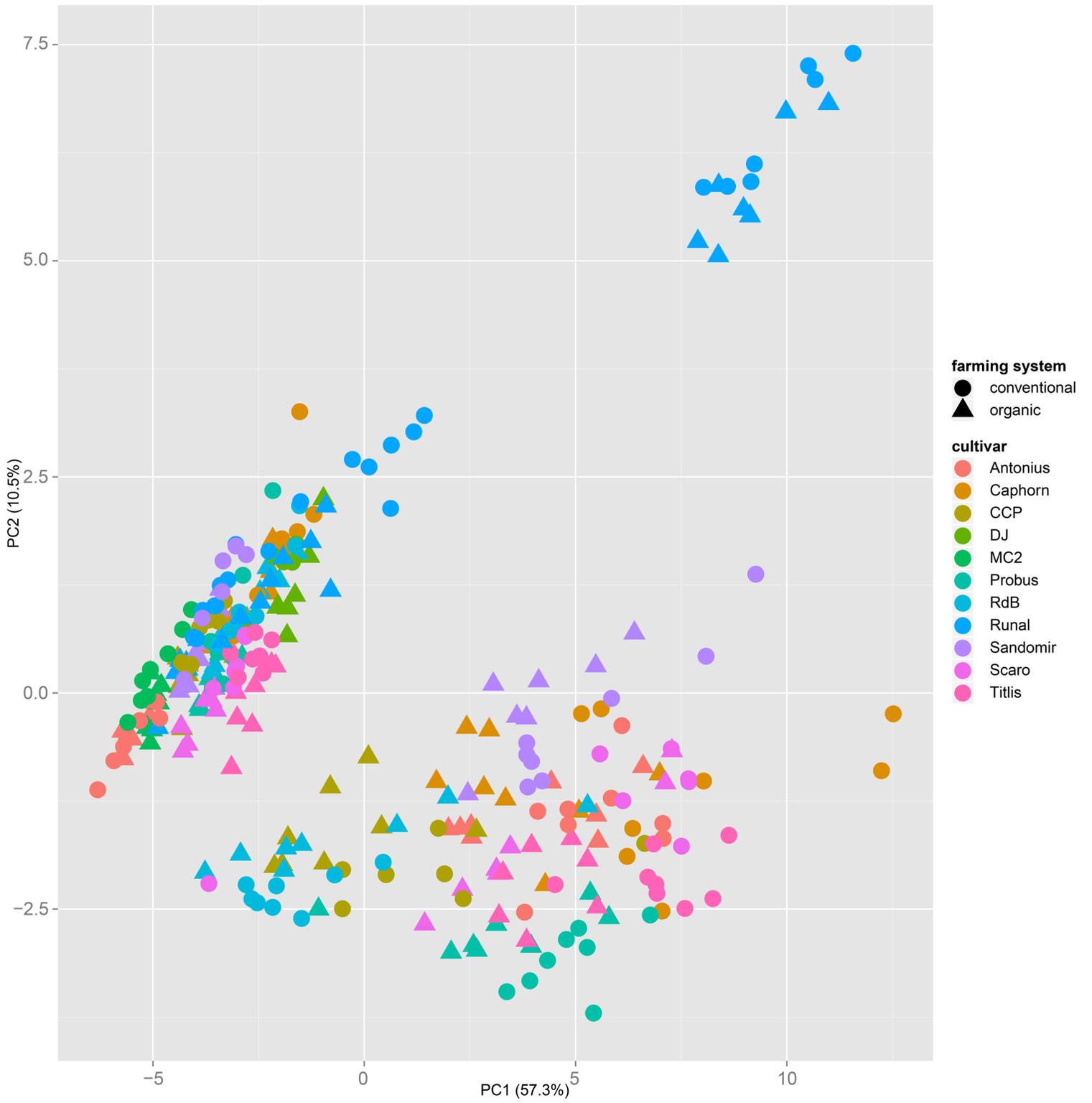


Figure 2.TIF

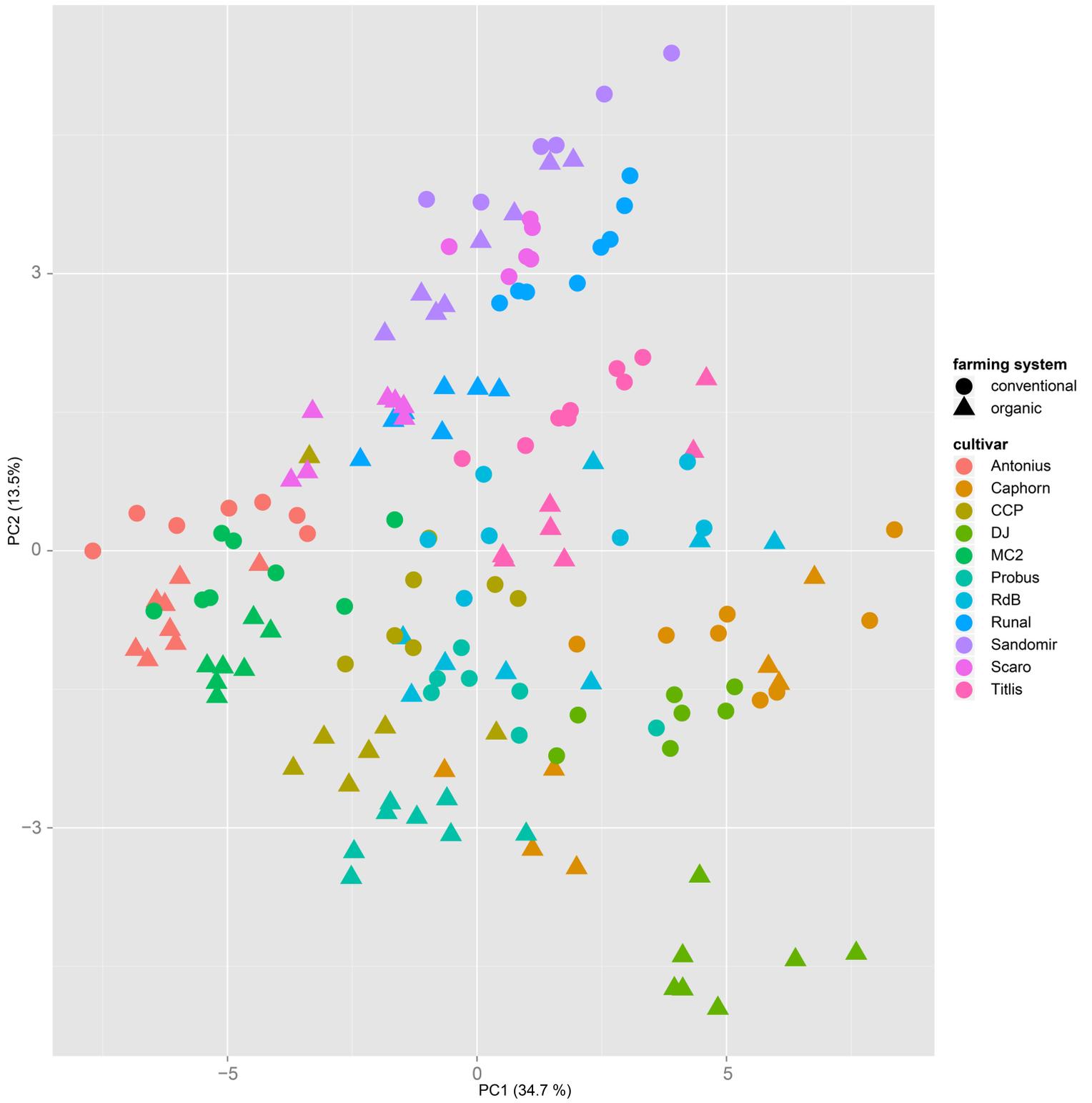


Figure 3.TIF

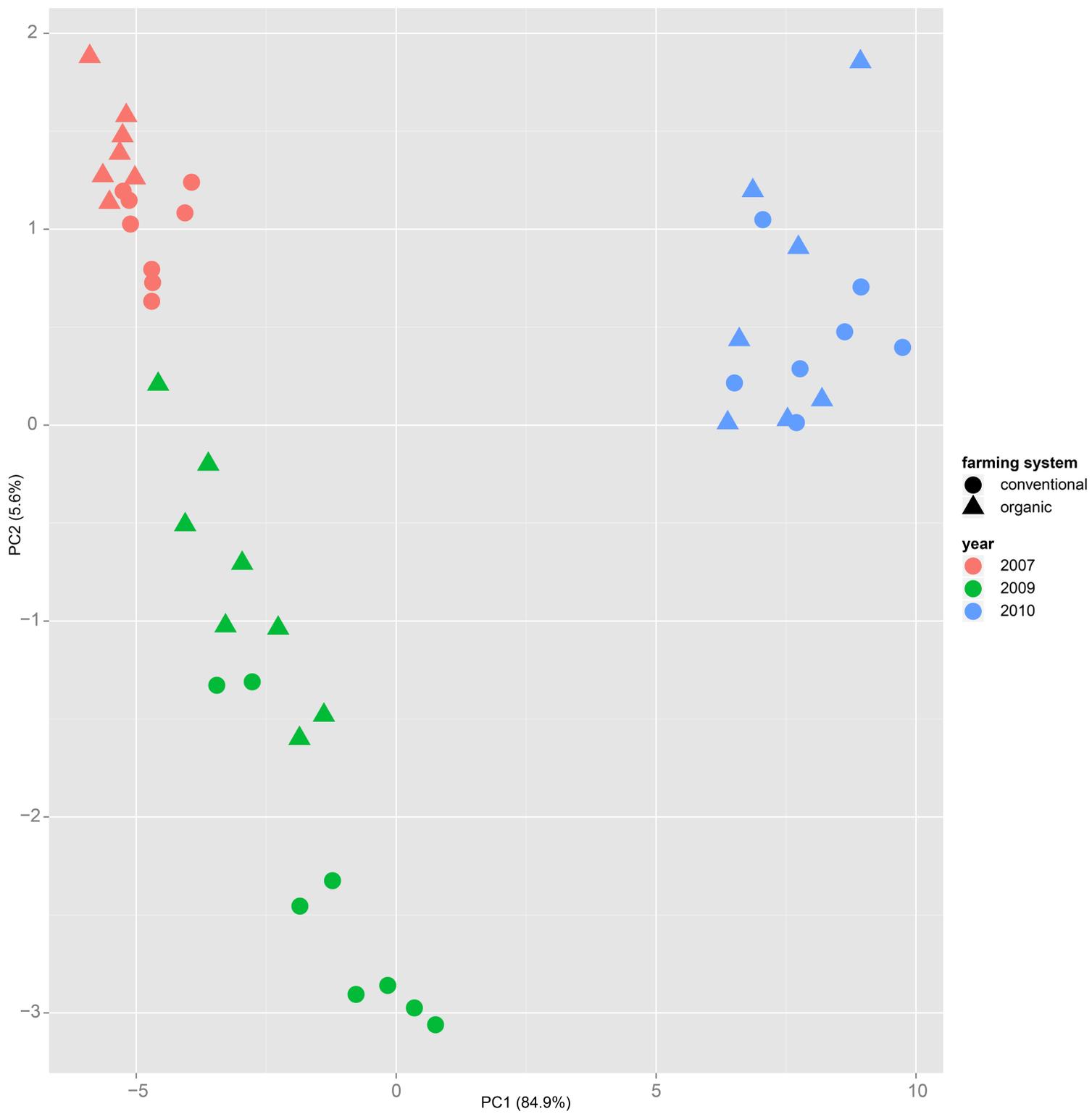


Figure 4.TIF

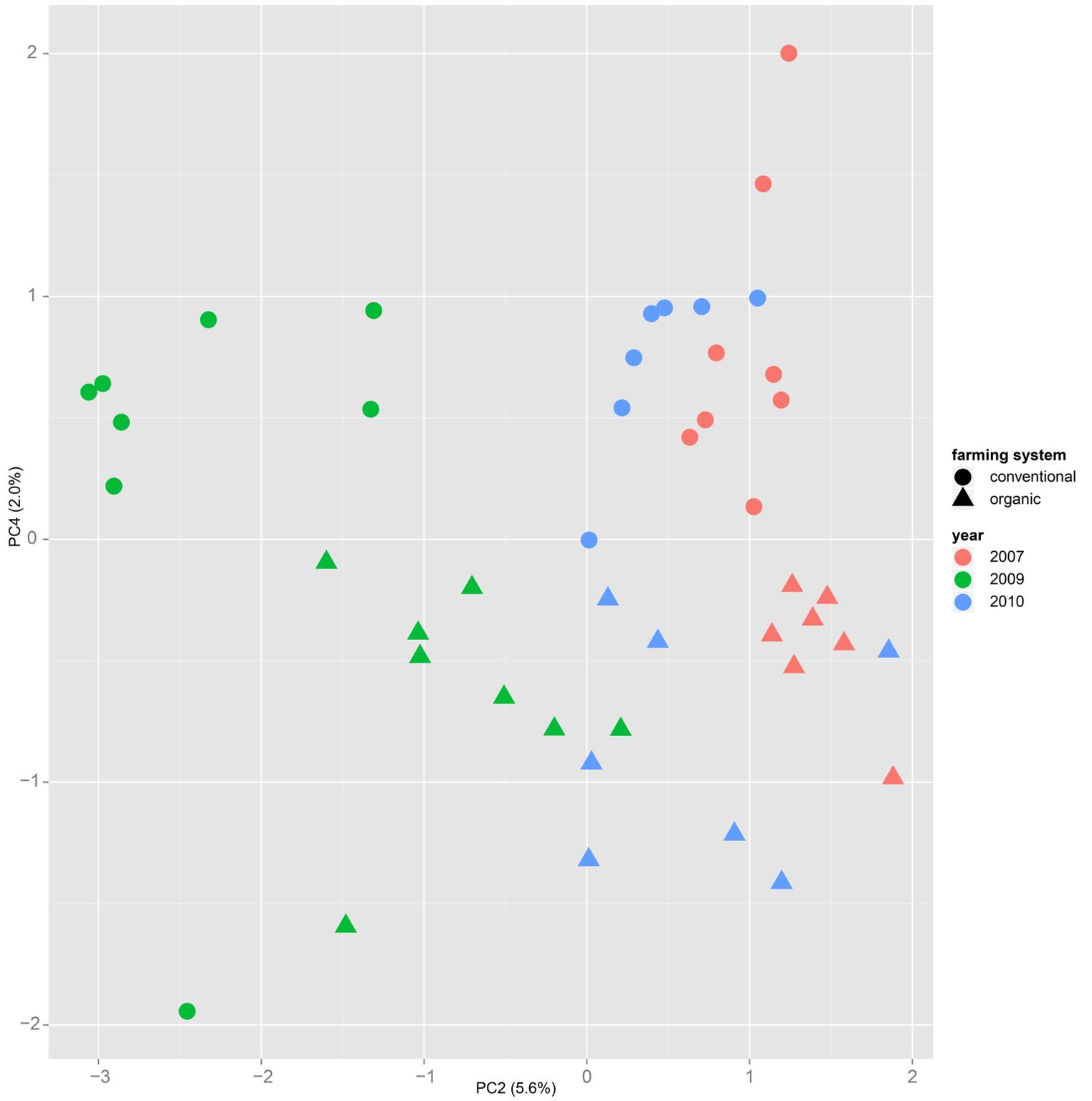


Figure 5.TIF

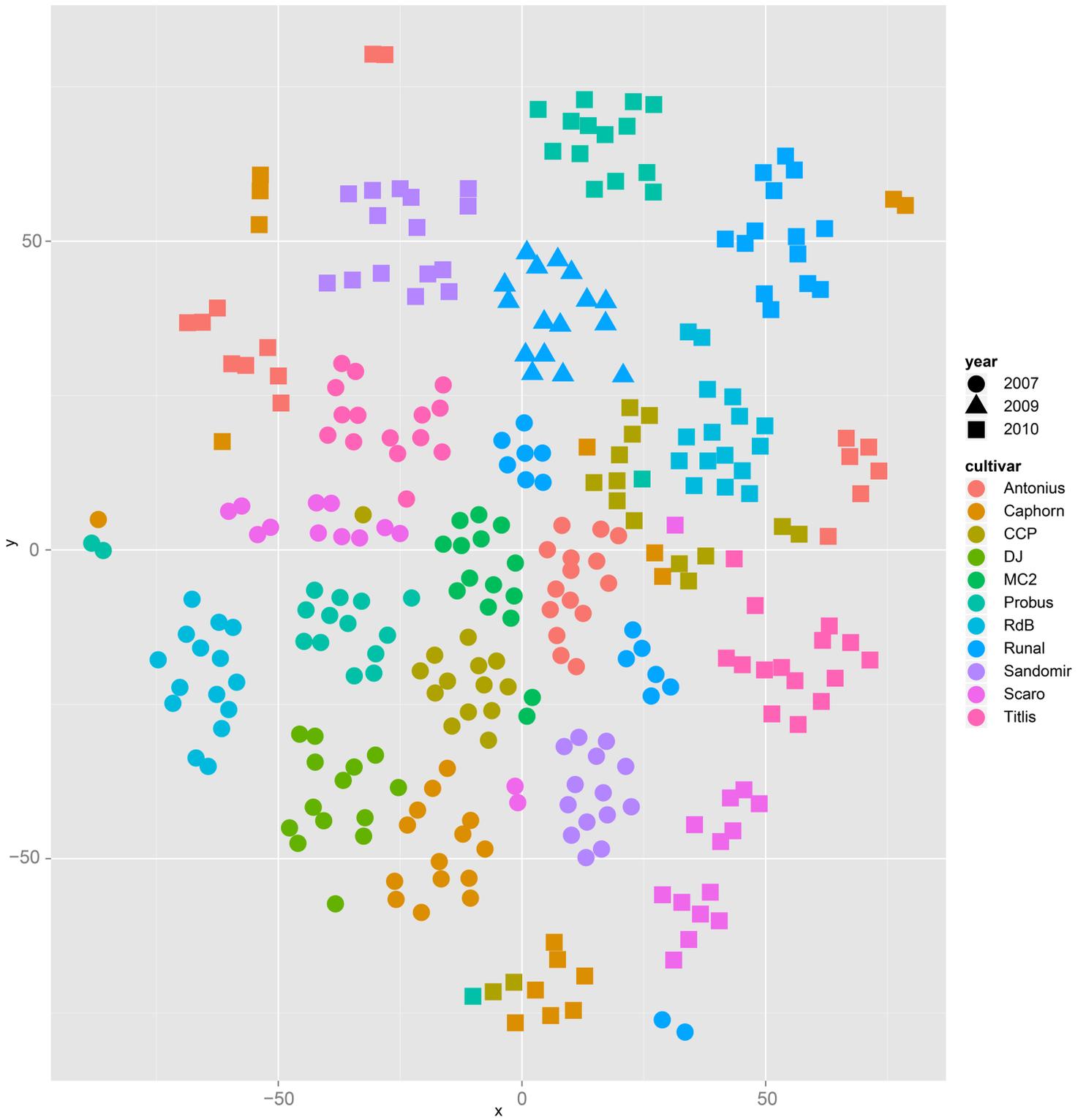


Figure 6.TIF

