

dpcR a Swiss-army knife for the analysis of digital PCR experiments

Michał Burdukiewicz^{1,5}, Jim Huggett², Alexandra Whale², Bart K.M. Jacobs³, Lieven Clement³, Piotr Sobczyk¹, Andrej-Nikolai Spiess⁴, Peter Schierack⁵, Stefan Rödiger^{5,*}

¹Department of Genomics, Faculty of Biotechnology, University of Wrocław, Wrocław, Poland and ²Molecular and Cell Biology Team, LGC, Teddington, United Kingdom and ³Department of Applied Mathematics, Computer Science and Statistics, Ghent University, Belgium and ⁴University Medical Center Hamburg-Eppendorf, Hamburg, Germany and ⁵Faculty of Natural Sciences, Brandenburg University of Technology Cottbus–Senftenberg, Großenhainer Str. 57, 01968, Senftenberg, Germany

Received January 1, 2009; Revised February 1, 2009; Accepted March 1, 2009

ABSTRACT

The digital Polymerase Chain Reaction (dPCR) is emerging in all research areas, such as life-sciences and diagnostics since it enables an absolute quantification of nucleic acids. Different approaches for statistical analysis were proposed. However, most analysis is done in closed source software as provided by the vendors. This makes it harder to compare results, such as the confidence interval estimates. An unified open software framework for reproducible research is not available.

To perform dPCR analysis we implemented peer-review statistical methods and plots into the dpcR framework, based on the sophisticated statistical computing environment R. dpcR is versatile open source cross-platform software framework, which provides functions to process dPCR data independent of the hardware. Our software can be used for data analysis and presentation, as framework for novel technical developments and as reference for statistical methods in dPCR analysis. Features such as functions to estimate the underlying Poisson process, calculation of confidence intervals based on single samples as well as on replicates, a novel Generalized Linear Model-based procedure to compare digital PCR experiments and a spatial randomness test for assessing plate effects have been integrated. We use a plug-in like architecture and abstraction layers to make the framework usable for droplets and (real-time) chamber based technologies.

dpcR is implemented with interfaces to the command-line, graphical user interfaces and interactive web application. Therefore, it can be used by novices in a graphical user interface or by experts via a command-line interface. The dpcR framework can be used to build a custom-made analyser according to the user requirements. dpcR is an

open framework, which can be easily adapted to the growing knowledge in dPCR.

INTRODUCTION

Real-time quantitative PCR (qPCR) is the standard approach to quantify nucleic acids (20). The quantification of the amplification is not done by determining a C_q-value derived from an amplification curve. qPCR is a well established and robust technology, which allows precise quantification of DNA material in high throughput fashion. However, the quantification by qPCR is challenging at very low and very high concentrations. At low DNA concentration Monte Carlo effect play a role and at high concentration inhibition processes start to dominate the qPCR. Thus, the qPCR is only usable in the working range of the calibrator. In addition, pre-processing and data analysis is affected by numerous adverse effects (27, 29).

The digital PCR (dPCR) is an important contender for precise nucleic acids quantifications. The chemical basis of the dPCR is similar to the qPCR, which includes master-mix preparation and thermal cycling of the sample. Though, approaches based on isothermal amplification were also developed (23). In contrast to qPCR, the amplification reaction does not take place in a single reaction chamber. Rather its a process of clonal amplification in small separate “partitions” (e.g., nl volume droplets of water oil emulsions, chambers on micro structured chips). The number of positive partition in relation to the number of total partitions. By applying Poisson statistics it is possible to determine the number of the starting material in given volume. Therefore, the dPCR does not require an external calibration (23, 28). Since approximately ten years the digital PCR (dPCR) is gaining momentum in the mainstream user-base. dPCR will likely have the same impact as qPCR in the nucleic acid methodology (11, 18, 23). There is an intensive research on dPCR platforms with the overall aim to make to technology broadly usable, cheap, robust and to enable high sample throughput.

A first proposal for digital PCR like approach and the use of the Poisson distribution to quantify the number of molecules in a sample was shown by Ruano *et al.* 1990 (PNAS) with the single molecule dilution (SMD) PCR (26). In 1999 Vogelstein *et al.* (PNAS) described the first true digital

*To whom correspondence should be addressed. Tel: +49 357385936; Fax: +49 357385801; Email: stefan.roediger@b.tu-cottbus.de

PCR (32). Application of the dPCR cover all applications of conventional qPCR, including investigation of alleles, gene expression analysis and absolute quantification of PCR products. For absolute quantification the qPCR relied on an external calibrator (calibration curve) which was derived serial decadic dilution (e.g., 1:10 → 1:100 → 1:1000) of a known target input quantity. The real-time monitoring of the PCR product formation enabled to determine quantification points (Cq). The Cq are strictly related to the input quantity. A simple arithmetic operation (after logarithmic transformation of the concentration) is sufficient to determine any nucleic acid quantity (10).

The dPCR has some principle assumptions and fundamental properties. First of all the chemical reaction should be not affected by inhibitors. The distribution of the single molecule target regions follows a Poisson distribution. The Poisson distribution appears like a normal distribution but without negative values and being zero the lowest. First a large number (n) of amplifications reactions as required to have a high statistical power. Therefore in practical terms a massive number of PCR reactions is needed. For Poisson distributions an n of XY (get reference from table/text book form statistics/biostatistics?) is considered large. Second that the molecules required for the amplification amplifications reactions are randomly distributed in the compartments. Visual analysis, Ripley's K functions or ??? can be used to test for randomness of the reaction and thus to exclude the clustering of of positive reactions. A clustering of positive wells might be due to sample loading or analysis process (systematical error). The outcome of an amplification can be no amplification at all (less than 1 copy per volume), an unsaturated reaction with a binary/"multinary" amplification (usable to calculate the "concentration") or a saturated reaction where virtually all compartments are positive.

Calculation of the "Concentration" Reference to "Supplement"

Calculation of the uncertainty To determine the uncertainty of the calculations two approach have been proposed in the peer-review literature (2, 7). The uncertainty is dependent on the number of PCR reactions (reference to *dpcR* functions). Reference to "Supplement" and *dpcR* functions.

Recently, Mathew *et al.* published the open access bioinformatic pipeline, designated **definetherain** (13). The tool is coded in JavaScript and has been made available for free in a web browser.

Interactive use and graphical representation with *shiny* (4).

Import and export of results figures and data.

There are currently two technical approaches to dPCR. dPCRs may use (microfluidic)chambers or emulsion based droplets (QX200TM(Bio-Rad), RainDropTMSystem (RainDance)). Chamber based dPCR systems have fixed geometries, including the volume of the reaction chambers. Despite the fact that dPCRs is an endpoint analysis the chamber based technologies allow generally the real-time monitoring of the amplification reaction and subsequent confirmation of the amplification reaction be melting curve analysis. Thus, such technologies enable easier trouble shooting and quality management of the data. However, the downside of these technologies is the fixed limited number of compartments and the price. The emulsion based dPCRs are

easier to perform since the compartments are generated by microfluidic technologies and have practically no limitation regarding the number of compartments. This results in a higher statistical power to quantify small differences in sample quantities. The emulsion chambers are made of water-in-oil emulsions with similar sizes.

There is a need for an vendor independent data analysis. For example, others have written custom made scripts for data analysis in **Mathematica** (Wolfram Research), **MS EXCEL** (Microsoft) or **R** (5, 6, 30, 31). However, this is of limited use, since the solutions are tied to a specific dPCR platform (e.g., droplet dPCR by Bio-Rad), operating system platform for data analysis and only usable for a single task. Moreover, we found no software packages with GUIs and bindings to a sophisticated statistical computing environment for reproducible research.

MATERIALS AND METHODS

Materials subsection one

We have chosen **R** because it is the *lingua franca* in biostatistics and broadly used in other disciplines. Since all software is open source it is possible to track numerical errors (23). There are many packages in existence which enable the fast development of new methods and plotting facilities. As most **R** packages depend on one or more other packages (19) depends *dpcR* on other packages, resulting in a complex network of recursive dependencies. Core packages *qpcR* (21), *shiny* (4), *MBmca* (22), *chipPCR* (24) and further packages as shown in the dependency graph (Supplement XYZ).

One basic design decision was to structure specific properties of digital PCR systems (droplet vs. chamber) in auxiliary functions and to perform central calculation specific to Poisson statistics in independent main functions. Chamber digital PCR systems fundamentally rely on the proper preprocessing of qPCR data. We have chosen to implement the core functionality by a dependency to the *qpcR* **R** package (21). The main functions (e.g., for analysis, simulations, plotting), several auxiliary helper functions (e.g., data import) and data set of different dPCR systems are listed in Table XY. Further dependencies to 3rd party packages include *pracma*, ... (see Figure 3). See the vignette for details.

The workflow is shown Figure 1.

The GUI employs advanced plots based on *ggplot2* (14).

Materials subsubsection one. qPCR dPCR Number of copies/DNA per volume (e.g., ng/l, copies/l) total number of compartments * ln (...)

Two-sided exact tests and matching confidence intervals for discrete data (8)

Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing (1)

Interval Estimation for a Binomial Proportion (3)

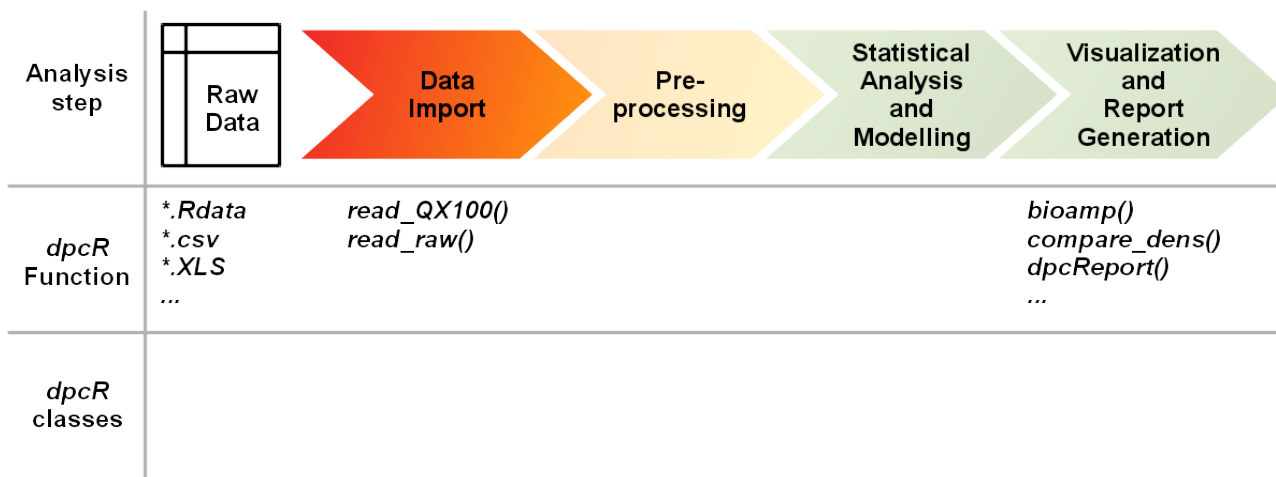


Figure 1. *dpcR* workflow. The diagram shows main functions available at each step of a dPCR data analysis.

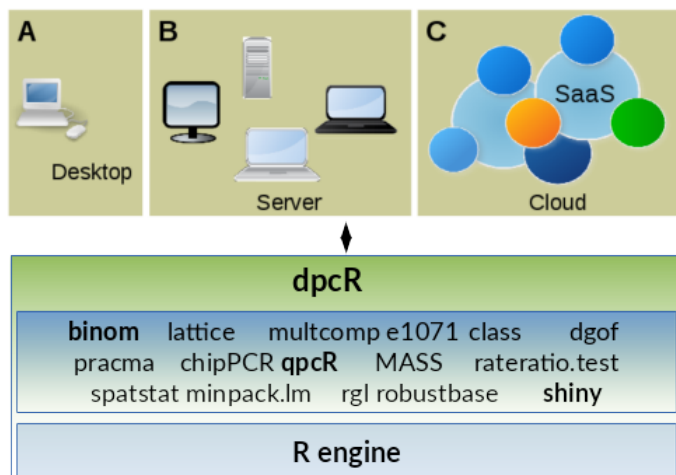


Figure 2. Modular software framework structure. *dpcR* is typically run from a desktop computer or a server. The software can be operated by an GUI/IDE application such as **RStudio** or **RKward**. The *dpcR* package has dependencies to other **R** packages (middle layer). The functionality shared between the packages enables repaid addition and expansion of functionality.

Integration in third party software

We aimed for a form factor (e.g., smart phone, tablet, desktop PCR) and operating system independent implementation of a graphical user interface. Moreover, there exit several GUI technologies for **R** (25). An interesting feature of the *shiny* technology is the automatic integration in environments, which support HTML5 and ECMAScript. This can be a modern web browser or an **R** IDE/GUI such as **RKward** (Figure 4) or **RStudio**.

Materials subsection Statistical power - Monte Carlo simulations

The proposed framework was evaluated in three Monte Carlo experiments (2000 times repetitions each) with accordingly 1000, 5000 (results not shown) and 10k partitions. During

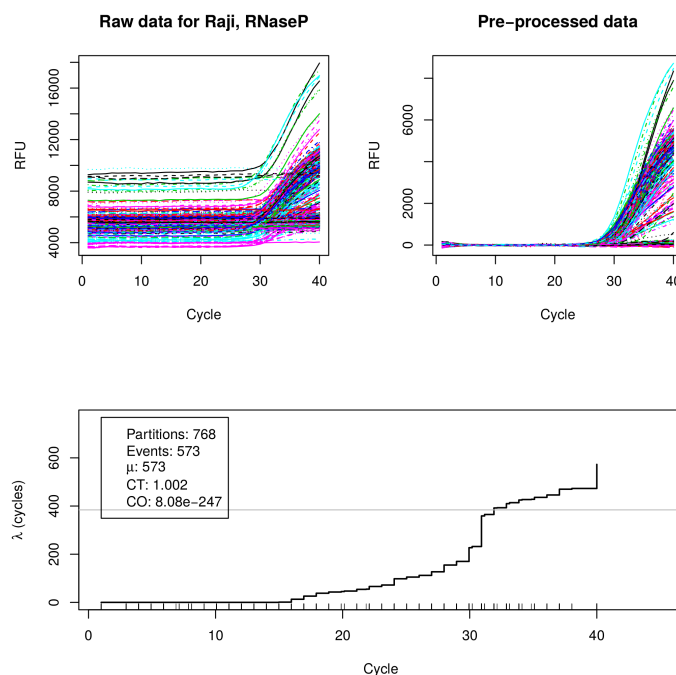


Figure 3. Uncover characteristics of dPCR data. Selected dPCR platforms are qPCR platforms at the same time. The function *qpcr2pp* uses the qPCR amplification curve data and interprets them as dPCR (Poisson process). A) Raw data of The function were B) preprocessed (baselined, smoothed) with functions from the *chipPCR* package and C) finally analyzed (Cq calculation → binarize) with the *qpcr2pp* (qPCR to Poisson process) function from the *dpcR* package.

each repetition of the Monte Carlo scheme, a set of partitions was randomly generated (7) with a determined number of molecules ('Base number of molecules' on X-axis). The set was copied and a number of molecules ('Added number of molecules' on Y-axis) was added to randomly chosen partitions. Two obtained arrays were compared using the

4 Nucleic Acids Research, 2009, Vol. 37, No. 12

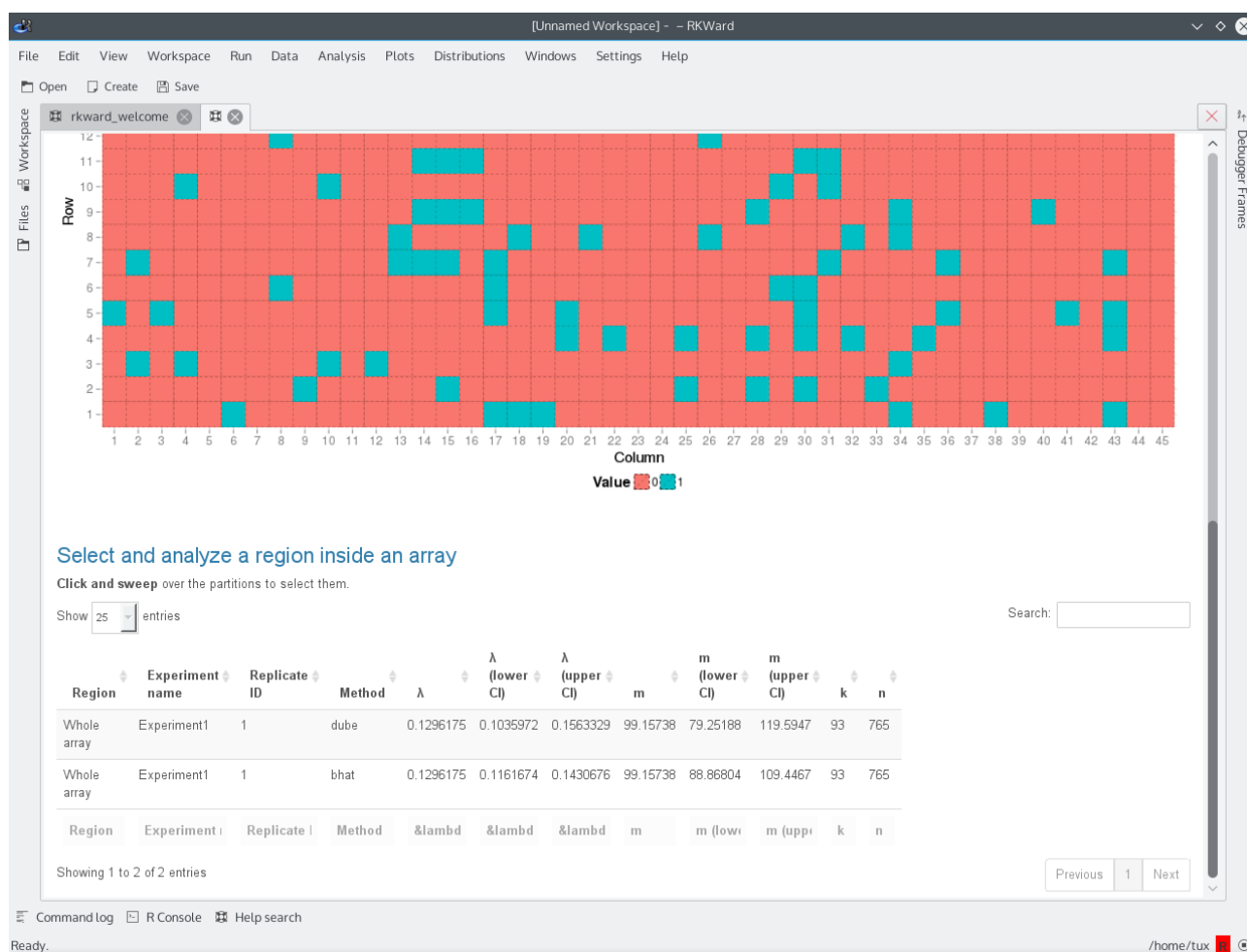


Figure 4. *dpcrReport()* function running in the graphical user interface and integrated development environment RKward.

proposed method. The mean p-values alongside with their standard deviation are presented in the chart below.

RESULTS

Reading or importing data

R has a rich set of tool to arrange data (reshape (35)) in order to prepare them for the analysis. This is important when it comes to the question how experiments should be treated. It is possible to analyze the PCR reaction the panels independently (effect on CI and uncertainty) or to pool/aggregate all reactions (effect on CI and uncertainty) to achieve higher sensitivity/certainty. (see Table 1).

De novo creation of dpcR objects from raw data

In experimental setups user have the need to transform their raw data in a processable format. Since the **R** environment is cross-platform an ubiquitously used we aimed to ease this creation. In particular, this is relevant for reproducible research. The *dpcR* framework covers the types typically used in laboratories.

Public data sets

dpcR includes data sets or refers to additional **R** packages for testing purposes. The data originate from different dPCR and qPCR systems and were either published previously (23, 24, 33, 34) or *de novo* generated.

Export of analysis results

Since *dpcR* is based on the **R** environment all facilities for a report generation are usable as described before (23). (see Table 1).

Table 1. Structured vendor export data formats handled by *dpcR* v. 0.5 and later.

Vendor	System	Format	Type
Bio-Rad	QX100 & QX200	CSV	Summary export
Fluidigm	BioMark	CSV	Summary export
Formulatrix	Constellation Digital PCR	CSV	Summary export

The number of structured export data formats handled by *dpcR* is growing. Numerous data formats can be processed with the functionality provided by the **R** environment (see (23)). CSV, comma separated values.

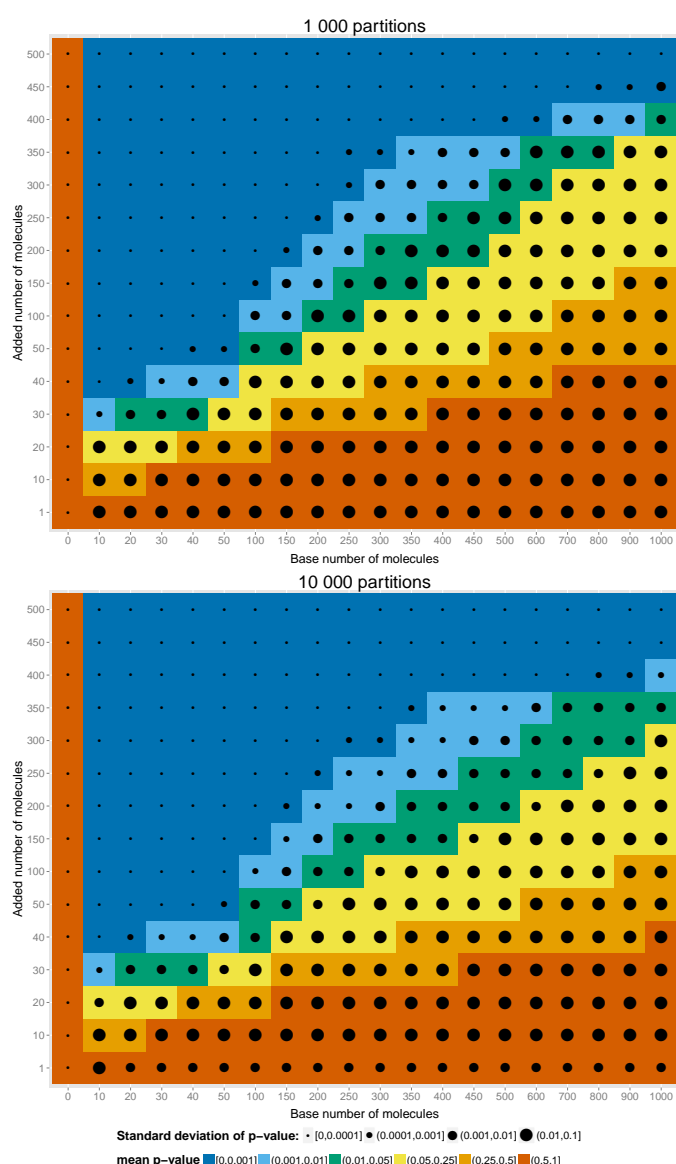


Figure 5. Our method, based on GLM, predicts estimated means copies per partitions using Poisson or binomial regression. Afterwards, estimates are compared against themselves using t-test. Obtained p-values and confidence intervals do not require further correction, because the familywise error is controlled through the whole analysis.

Availability

The *dpcR* framework is available as open source software package as part of the Bioconductor project (9). The source code is open source (GPL-3 or later). The stable version is hosted at <http://cran.r-project.org/web/packages/dpcR> and the source code is available from <https://github.com/michbur/dpcR>.

Documentation

All functions of the *dpcR* package have its own documentation package, which specifies the input types, classes, parameters

and output formats. The documentation is available as standard **R** package reference manual and as vignette.

DISCUSSION

Currently, there exist different dPCR analysis software solutions provided by the vendors. But most of the software packages are designed black boxes, which prevent deep insight into the data processing step. Other and we think that scientific software should be open (12, 17, 23). In addition, most of the software solutions are aimed to be used in very specific scenarios and a mutual exclusive to alternative platforms (e.g., droplet vs. chamber-based). We have chosen **R** because it is the *lingua franca* in biostatistics and broadly used in other disciplines (23). We developed the *dpcR* package, which is a software framework for analysis of dPCR. *dpcR* provides the scientific community a broadly applicable tool for teaching purposes, data analysis and theoretical research based on simulations. Our software framework can be used to accelerate the development of new approaches to dPCR.

We implemented numerous statistical methods for dPCR and suggest the introduction of a standardized nomenclature for dPCR. The package enables the simulations and predictions of dPCR reactions and the analysis of previously run dPCRs.

Functions included may be used to simulate dPCRs, perform statistical data analysis, plotting of the results and simple report generation.

We decided not to implement algorithms for clustering and “rain” (positive droplets) definition of droplet dPCR data. This is because, there are several **R** packages from flow-cytometer research. Implementations range from manual to automatic clustering (15, 16, 31). Moreover, discussion with our peers and the literature suggest that a consensus of an appropriate method for dPCR is not available.

CONCLUSION

In conclusion, *dpcR* provides means to understand how digital PCR works, to design, simulate and analyze experiments, and to verify their results (e.g., confidence interval estimation), which should ultimately improve reproducibility. We have built what we believe to be the first unified, cross-platform, dMIQE compliant, open source software framework for analyzing dPCR experiments. Our *dpcR* framework is targeted at a broad user base including end users in clinics, academics, developers, and educators. We implemented existing statistical methods for dPCR and suggest the introduction of a standardized dPCR nomenclature. Our framework is suitable for teaching and includes references for an elaborated set of methods for dPCR statistics. Our software can be used for (I) data analysis and visualization in research, (II) as software framework for novel technical developments, (III) as platform for teaching this new technology and (IV) as reference for statistical methods with a standardized nomenclature for dPCR experiments. The framework enables the simulations and predictions of Poisson distribution for dPCR scenarios, the analysis of previously run dPCRs. Due to the plug-in structure of the software it is possible to build custom-made analyzers.

Our open framework includes to invitation to the scientific community to join and support the development of *dpcR*.

ACKNOWLEDGEMENTS

Grateful thanks belong to the **R** community and the **RStudio** developers.

FUNDING

This work was funded by the Federal Ministry of Education and Research (BMBF) InnoProfile–Transfer–Projekt 03 IPT 611X.

Conflict of interest statement. None declared.

REFERENCES

1. Yoav Benjamini and Yosef Hochberg. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):289–300, 1995.
2. Somanath Bhat, Jan Herrmann, Paul Armishaw, Philippe Corbisier, and Kerry R Emslie. Single molecule detection in nanofluidic digital array enables accurate measurement of DNA copy number. *Analytical and Bioanalytical Chemistry*, 394(2):457–467, May 2009.
3. Lawrence D. Brown, T. Tony Cai, and Anirban DasGupta. Interval estimation for a binomial proportion. *Statist. Sci.*, 16(2):101–133, May 2001.
4. Winston Chang, Joe Cheng, JJ Allaire, Yihui Xie, and Jonathan McPherson. *shiny: Web Application Framework for R*, 2015. R package version 0.12.2.
5. David Dobnik, Björn Spilsberg, Alexandra Bogoalec Koir, Arne Holst-Jensen, and Jana el. Multiplex Quantification of 12 European Union Authorized Genetically Modified Maize Lines with Droplet Digital Polymerase Chain Reaction. *Analytical Chemistry*, 87(16):8218–8226, August 2015.
6. Tanja Dreö, Manca Pirc, Iva Ramak, Jernej Pavi, Mojca Milavec, Jana el, and Kristina Gruden. Optimising droplet digital PCR analysis approaches for detection and quantification of bacteria: a case study of fire blight and potato brown rot. *Analytical and Bioanalytical Chemistry*, 406(26):6513–6528, August 2014.
7. Simant Dube, Jian Qin, and Ramesh Ramakrishnan. Mathematical analysis of copy number variation in a DNA sample using digital PCR on a nanofluidic device. *PloS one*, 3(8):e2876, 2008.
8. Michael Fay. Two-sided exact tests and matching confidence intervals for discrete data. *Proceedings of the National Academy of Sciences of the United States of America*, 2(1):53–58, June 2010.
9. Robert C Gentleman, Vincent J Carey, Douglas M Bates, Ben Bolstad, Marcel Dettling, Sandrine Dudoit, Byron Ellis, Laurent Gautier, Yongchao Ge, Jeff Gentry, Kurt Hornik, Torsten Hothorn, Wolfgang Huber, Stefano Iacus, Rafael Irizarry, Friedrich Leisch, Cheng Li, Martin Maechler, Anthony J Rossini, Gunther Sawitzki, Colin Smith, Gordon Smyth, Luke Tierney, Jean Y H Yang, and Jianhua Zhang. Bioconductor: open software development for computational biology and bioinformatics. *Genome biology*, 5(10):R80, 2004.
10. Jim F. Huggett, Simon Cowen, and Carole A. Foy. Considerations for Digital PCR as an Accurate Molecular Diagnostic Tool. *Clinical Chemistry*, page clinchem.2014.221366, October 2014.
11. Jim F. Huggett, Justin OGrady, and Stephen Bustin. qPCR, dPCR, NGS A journey. *Biomolecular Detection and Quantification*, 3:A1–A5, March 2015.
12. Darrel C. Ince, Leslie Hatton, and John Graham-Cumming. The case for open computer programs. *Nature*, 482(7386):485–488, February 2012.
13. Mathew Jones, James Williams, Kathleen Gärtner, Rodney Phillips, Jacob Hurst, and John Frater. Low copy target detection by Droplet Digital PCR through application of a novel open access bioinformatic pipeline, definetherain. *Journal of Virological Methods*, 202(100):46–53, June 2014.
14. David Kahle and Hadley Wickham. ggmap: Spatial visualization with ggplot2. *The R Journal*, 5(1):144–162, 2013.
15. Nolwenn Le Meur. Computational methods for evaluation of cell-based data assessmentBioconductor. *Current Opinion in Biotechnology*, 24(1):105–111, February 2013.
16. Mehrnoush Malek, Mohammad Jafar Taghiyar, Lauren Chong, Greg Finak, Raphael Gottardo, and Ryan R. Brinkman. flowdensity: reproducing manual gating of flow cytometry data by automated density-based cell population identification. *Bioinformatics*, 31(4):606–607, 2015.
17. A. Morin, J. Urban, P. D. Adams, I. Foster, A. Sali, D. Baker, and P. Sliz. Shining Light into Black Boxes. *Science*, 336(6078):159–160, April 2012.
18. Alexander A. Morley. Digital PCR: A brief history. *Biomolecular Detection and Quantification*, 1(1):1–2, 2014.
19. Jeroen Ooms. Directions for improved dependency versioning in R. *The R Journal*, 5(1):197–207, 2013.
20. Stephan Pabinger, Stefan Rödiger, Albert Kriegner, Klemens Vierlinger, and Andreas Weinhäusel. A survey of tools for the analysis of quantitative PCR (qPCR) data. *Biomolecular Detection and Quantification*, 1(1):23–33, 2014.

21. Christian Ritz and Andrej-Nikolai Spiess. qpcR: an R package for sigmoidal model selection in quantitative real-time polymerase chain reaction analysis. *Bioinformatics*, 24(13):1549–1551, January 2008.
22. Stefan Rödiger, Alexander Böhm, and Ingolf Schimke. Surface Melting Curve Analysis with R. *The R Journal*, 5(2):37–53, December 2013.
23. Stefan Rödiger, Michał Burdukiewicz, Konstantin A. Blagodatskikh, and Peter Schierack. R as an Environment for the Reproducible Analysis of DNA Amplification Experiments. *The R Journal*, 7(2):127–150, 2015.
24. Stefan Rödiger, Michał Burdukiewicz, and Peter Schierack. chipPCR: an R package to pre-process raw data of amplification curves. *Bioinformatics*, 31(17):2900–2902, 2015.
25. Stefan Rödiger, Thomas Friedrichsmeier, Prasenjit Kapat, and Meik Michalke. RKWard: A Comprehensive Graphical User Interface and Integrated Development Environment for Statistical Analysis with R. *Journal of Statistical Software*, 49(9):1–34, 2012.
26. G. Ruano, K. K. Kidd, and J. C. Stephens. Haplotype of multiple polymorphisms resolved by enzymatic amplification of single DNA molecules. *Proceedings of the National Academy of Sciences*, 87(16):6296–6300, January 1990.
27. Jan M Ruijter, Michael W Pfaffl, Sheng Zhao, Andrej N Spiess, Gregory Boggy, Jochen Blom, Robert G Rutledge, Davide Sisti, Antoon Lievens, Katleen De Preter, Stefaan Derveaux, Jan Helleman, and Jo Vandesompele. Evaluation of qPCR curve analysis methods for reliable biomarker discovery: bias, resolution, precision, and implications. *Methods (San Diego, Calif.)*, 59(1):32–46, 2013.
28. David A. Selck, Mikhail A. Karymov, Bing Sun, and Rustem F. Ismagilov. Increased Robustness of Single-Molecule Counting with Microfluidics, Digital Isothermal Amplification, and a Mobile Phone versus Real-Time Kinetic Measurements. *Analytical Chemistry*, 85(22):11129–11136, November 2013.
29. Andrej-Nikolai Spiess, Claudia Deutschmann, Michał Burdukiewicz, Ralf Himmelreich, Katharina Klat, Peter Schierack, and Stefan Rödiger. Impact of Smoothing on Parameter Estimation in Quantitative DNA Amplification Experiments. *Clinical Chemistry*, 61(2):379–388, January 2015.
30. Matthew C Strain, Steven M Lada, Tiffany Luong, Steffney E Rought, Sara Gianella, Valeri H Terry, Celsa A Spina, Christopher H Woelk, and Douglas D Richman. Highly precise measurement of HIV DNA by droplet digital PCR. *PloS one*, 8(4):e55943, 2013.
31. Wim Trypsteen, Matthijs Vynck, Jan De Neve, Pawel Bonczkowski, Maja Kiselinova, Eva Malatinkova, Karen Vervisch, Olivier Thas, Linos Vandekerckhove, and Ward De Spiegelaere. ddpcRquant: threshold determination for single channel droplet digital PCR experiments. *Analytical and Bioanalytical Chemistry*, 407(19):5827–5834, July 2015.
32. B Vogelstein and K W Kinzler. Digital PCR. *Proceedings of the National Academy of Sciences of the United States of America*, 96(16):9236–9241, August 1999.
33. Alexandra S Whale, Jim F Huggett, Simon Cowen, Valerie Speirs, Jacqui Shaw, Stephen Ellison, Carole A Foy, and Daniel J Scott. Comparison of microfluidic digital PCR and conventional quantitative PCR for measuring copy number variation. *Nucleic Acids Research*, 40(11):e82, June 2012.
34. Richard A White, 3rd, Paul C Blainey, H Christina Fan, and Stephen R Quake. Digital PCR provides sensitive and absolute calibration for high throughput sequencing. *BMC genomics*, 10:116, 2009.
35. Wickham and Hadley. Reshaping data with the reshape package. *Journal of Statistical Software*, 21(12), 2007.