

A peer-reviewed version of this preprint was published in PeerJ on 16 March 2018.

[View the peer-reviewed version](https://peerj.com/articles/4473) (peerj.com/articles/4473), which is the preferred citable publication unless you specifically need to cite this preprint.

Ahmed M, Kim DR. 2018. pcr: an R package for quality assessment, analysis and testing of qPCR data. PeerJ 6:e4473
<https://doi.org/10.7717/peerj.4473>

pcr: an R package for quality assessment, analysis and testing of qPCR data

Mahmoud Ahmed¹, Deok Ryong Kim^{Corresp. 2}

¹ Department of Biochemistry and Convergence Medical Sciences, Gyeongsang National University School of Medicine, Jinju, GyeongNam, South Korea

² Department of Biochemistry and Convergence Medical Sciences, Institute of Health Sciences, Gyeongsang National University School of Medicine, Jinju, GyeongNam, South Korea

Corresponding Author: Deok Ryong Kim
Email address: drkim@gnu.ac.kr

Background. Real-time quantitative PCR (qPCR) is a broadly used technique in the biomedical research. Currently, few different analysis models are used to determine the quality of data and to quantify the mRNA level across the experimental conditions.

Methods. We developed an R package to implement methods for quality assessment, analysis and testing qPCR data for statistical significance. Double Delta *CT* and standard curve models were implemented to quantify the relative expression of target genes from *CT* in standard qPCR control-group experiments. In addition, calculation of amplification efficiency and curves from serial dilution qPCR experiments are used to assess the quality of the data. Finally, two-group testing and linear models were used to test for significance of the difference in expression control groups and conditions of interest. **Results.** Using two datasets from qPCR experiments, we applied different quality assessment, analysis and statistical testing in the pcr package and compared the results to the original published articles. The final relative expression values from the different models, as well as the intermediary outputs, were checked against the expected results in the original papers and were found to be accurate and reliable. **Conclusion.** The pcr package provides an intuitive and unified interface for its main functions to allow biologist to perform all necessary steps of qPCR analysis and produce graphs in a uniform way.

1 pcr: an R package for quality assessment, 2 analysis and testing of qPCR data

3 Mahmoud Ahmed¹ and Deok Ryong Kim²

4 ^{1,2}Department of Biochemistry and Convergence Medical Sciences and Institute of
5 Health Sciences, Gyeongsang National University School of Medicine, Jinju, Republic
6 of Korea 527-27

7 Corresponding author:

8 Deok Ryong Kim²

9 Email address: drkim@gnu.ac.kr

10 ABSTRACT

11 **Background.** Real-time quantitative PCR (qPCR) is a broadly used technique in the biomedical research.
12 Currently, few different analysis models are used to determine the quality of data and to quantify the
13 mRNA level across the experimental conditions. **Methods.** We developed an R package to implement
14 methods for quality assessment, analysis and testing qPCR data for statistical significance. Double Delta
15 *CT* and standard curve models were implemented to quantify the relative expression of target genes
16 from *CT* in standard qPCR control-group experiments. In addition, calculation of amplification efficiency
17 and curves from serial dilution qPCR experiments are used to assess the quality of the data. Finally,
18 two-group testing and linear models were used to test for significance of the difference in expression
19 control groups and conditions of interest. **Results.** Using two datasets from qPCR experiments, we
20 applied different quality assessment, analysis and statistical testing in the pcr package and compared the
21 results to the original published articles. The final relative expression values from the different models, as
22 well as the intermediary outputs, were checked against the expected results in the original papers and
23 were found to be accurate and reliable. **Conclusion.** The pcr package provides an intuitive and unified
24 interface for its main functions to allow biologist to perform all necessary steps of qPCR analysis and
25 produce graphs in a uniform way.

26 INTRODUCTION

27 Real-time quantitative PCR (qPCR) is a commonly used technique to analyze the relative gene expression
28 level in the biomedical research. In most cases, small scale experiments are designed to quantify the
29 level of mRNA among experimental conditions. Some advanced machines and optimized protocols have
30 simplified the experiments to a highly efficient one-step process, allowing the effective analysis of a
31 large scale of qPCR data. However, all processes for assessing the quality of the data, performing the
32 analysis and reporting the results are not done in the most uniform way across the literature (Bustin
33 and Nolan, 2004). Different analysis models have been proposed and implemented in different software
34 environments (Pabinger et al., 2014). Furthermore, requirements and guidelines for reporting qPCR data
35 were independently introduced (Bustin et al., 2009).

36 In this report, we introduce an open source R package for performing quality assessment, modeling
37 and testing for statistical significance of qPCR data in a uniform way. In its current version, the pcr
38 package implement two methods for relative quantification of mRNA expression proposed originally by
39 Livak and Schmittgen (2001), in addition to the necessary steps to check the assumption of these methods.
40 Also, we implement a number of methods to check for statistical significance in qPCR data which were
41 introduced in SAS by Yuan et al. (2006). Finally, the package provides unified interface to make the
42 analysis accessible and the ability to make graphs of the different analysis steps for visual inspection and
43 preparation of publication-level figures. We start by describing the process for generating the data in the
44 original papers, briefly introduce the methods and apply them to the original data using the pcr.

45 MATERIALS & METHODS

46 Data Sources

47 To illustrate the usage of the pcr package and to apply it to qPCR data, we used real qPCR datasets from
 48 two published papers. In addition, we compared the results obtained by the pcr package to that of the
 49 original paper to ensure the reliability. At the first paper, Livak and Schmittgen (2001) obtained total RNA
 50 from human tissues; brain and kidney. c-myc and GAPDH primers were then used for cDNA synthesis
 51 and used as input in the PCR reaction. Seven different dilutions were used as input to the PCR reaction
 52 (three replicates each), this dataset was referred to as ct3 and shown in Table 1. Six replicates for each
 53 tissue were run in separate tubes. This dataset was referred to as ct1 through this document and shown
 54 along with the difference calculations in Table 2 and 3. At the second work, Yuan et al. (2006) extracted
 55 total RNA from *Arabidopsis thaliana* plant treated and control samples (24 samples each), and performed
 56 qPCR analyses using MT7 and ubiquitin primers. This dataset was referred to as ct4 and shown the
 57 results of the different testing methods that applied in the original paper in Table 4.

58 Statistical methods

59 In contrast with the absolute quantification of the amount of mRNA in a sample, the relative quantification
 60 uses a internal control (reference gene) and/or a control group (reference group) to quantify the mRNA of
 61 interest relative to these references. This relative quantification was sufficient to draw conclusions in most
 62 of the biomedical applications involving qPCR. A few methods were developed to perform these relative
 63 quantification. These methods generally require different assumptions and models for the analysis. The
 64 most common two of these methods were described here in the following sections.

65 *The comparative C_T methods*

66 The comparative C_T methods assume that the cDNA templates of the gene/s of interest as well as the
 67 control/reference gene have similar amplification efficiency, and also that this amplification efficiency is
 68 near perfect. Meaning, at a certain threshold during the linear portion of the PCR reaction, the amount
 69 of the gene of the interest and the control double each cycle. Another assumption is that the expression
 70 difference between two genes or two samples can be captured by subtracting one (gene or sample of
 71 interest) from another (reference). The final assumption is that the reference doesn't change with the
 72 treatment or the course in question. The formal derivation of the double delta C_T model is described here.
 73 Briefly, the $\Delta\Delta C_T$ is given by:

$$\Delta\Delta C_T = \Delta C_{T,q} - \Delta C_{T,cb} \quad (1)$$

74 And the relative expression by:

$$2^{-\Delta\Delta C_T} \quad (2)$$

75 Where $\Delta C_{T,q}$ is the difference in the C_T (or their average) of a gene of interest and a reference gene
 76 in a group of interest. $\Delta C_{T,cb}$ is the the difference in the C_T (or their average) of a gene of interest and a
 77 reference gene in a reference group. The error term is given by:

$$s = \sqrt{s_1^2 + s_2^2} \quad (3)$$

78 Where s_1 is the standard deviation of a gene of interest and s_2 is the standard deviation of a reference
 79 gene.

80 *Standard curve*

81 In comparison, this model doesn't assume perfect amplification but rather actively use the amplification
 82 in calculating the relative expression. So when the amplification efficiency of all genes are 100% both
 83 methods should give similar results. The standard curve method is applied using two steps. First, serial
 84 dilutions of the mRNAs from the samples of interest are used as input to the PCR reaction. The linear
 85 trend of the log input amount and the resulting C_T values for each gene are used to calculate an intercept
 86 and a slope. Secondly, these intercepts and slopes are used to calculate the amounts of mRNA of the

87 genes of interest and the control/reference in the samples of interest and the control sample/reference.
 88 These amounts are finally used to calculate the relative expression in a manner similar to the later method,
 89 just using division instead of subtraction. The formal derivation of the model is described here (Yuan
 90 et al., 2006). Briefly, The amount of RNA in a sample is given by:

$$\log amount = \frac{C_T - b}{m} \quad (4)$$

91 And the relative expression is given by:

$$10^{\log amount} \quad (5)$$

92 Where C_T is the cycle threshold of a gene. b is the intercept of $C_T - \log_{10}$ input amount. m is the
 93 slope of C_T . And the error term is given by:

$$s = (cv)(\bar{X}) \quad (6)$$

94 Where:

$$cv = \sqrt{cv_1^2 + cv_2^2} \quad (7)$$

95 Where s is the standard deviation. \bar{X} is the average. cv is the coefficient of variation or relative
 96 standard deviation.

97 **Statistical significance tests**

98 Assuming that the assumptions of the first methods are holding true, the simple t-test can be used to test
 99 the significance of the difference between two conditions (ΔC_T). t-test assumes, in addition, that the input
 100 C_T values are normally distributed and the variance between conditions are comparable. Wilcoxon test
 101 can be used if sample size is small, and those two last assumptions are hard to achieve.

102 Two use the linear regression here. A null hypothesis is formulated as following,

$$C_{T,target,treatment} - C_{T,control,treatment} = C_{T,target,control} - C_{T,control,control} \quad (8)$$

103 This is exactly the $\Delta\Delta C_T$ value as explained earlier. So the $\Delta\Delta C_T$ is estimated and the null is rejected
 104 when $\Delta\Delta C_T \neq 0$.

105 **Quality Assessment**

106 Fortunately, regardless of the method used in the analysis of qPCR data, The quality assessment can be
 107 done in a similar way. It requires an experiment similar to that of calculating the standard curve. Serial
 108 dilutions of the genes of interest and controls are used as input to the reaction and different calculations
 109 are made. The amplification efficiency is approximated by the linear trend between the difference between
 110 the C_T value of a gene of interest and a control/reference (ΔC_T) and the log input amount. This piece of
 111 information is required when using the $\Delta\Delta C_T$ model. Typically, the slope of the curve should be very
 112 small and the R^2 value should be very close to one. A value of the amplification efficiency itself is given
 113 by $10^{-1/slope}$ and should be close to 2. Other analysis methods are recommended when this is not the
 114 case. Similar curves are required for each gene using the C_T value instead of the difference for applying
 115 the standard curve method. In this case, a separate slope and intercept for each gene are required for the
 116 calculation of the relative expression.

117 **RESULTS & DISCUSSION**

118 **Availability & Installation**

119 The pcr packages is available on CRAN, the main repository for R packages and can be installed by
 120 invoking `install.packages` in an R ($\geq 3.4.2$) session. The package's source code is also available on
 121 github, <https://github.com/MahShaaban/pcr> along with the development version.

```

122 # install the pcr package from CRAN
123 install.packages('pcr')

```

124 The examples shown in this article are explained in greater details in the package vignette that can
 125 be accessed through `browseVignette('pcr')`. Moreover, the package documentation provides detailed
 126 instruction on the input and the output of each function (e.g. `?pcr_analyze`).

127 **Functionality & user interface**

128 The pcr package provides different methods for performing quality assessment, modeling and testing
 129 real-time qualitative PCR data through the unified interface of three functions `pcr_assess`, `pcr_analyze`
 130 and `pcr_test`, respectively.

131 **Quality Assessment**

132 `pcr_assess` provides two methods for assessing the quality of qPCR data. These are 'efficiency' and
 133 'standard_curve' to calculate the amplification efficiency and gene standard curves as described in the
 134 methods section. The following code block applies both methods to the dataset `ct3`, shown in Table 1.
 135 Using the argument `plot` as TRUE in the `pcr_assess` function provides the a graphic presentation of the
 136 amplification and the standard curves as shown in Figure 1.

```

137 # load required libraries
138 library(pcr)
139 library(ggplot2)
140 library(cowplot)
141 library(dplyr)
142 library(xtable)
143 library(readr)

144 # pcr_assess
145 ## locate and read data
146 fl <- system.file('extdata', 'ct3.csv', package = 'pcr')
147 ct3 <- read_csv(fl)
148
149 ## make a vector of RNA amounts
150 amount <- rep(c(1, .5, .2, .1, .05, .02, .01), each = 3)
151
152 ## calculate amplification efficiency
153 res1 <- pcr_assess(ct3,
154                   amount = amount,
155                   reference_gene = 'GAPDH',
156                   method = 'efficiency')
157
158 ## calculate standard curves
159 res2 <- pcr_assess(ct3,
160                   amount = amount,
161                   method = 'standard_curve')
162
163 ## retain curve information
164 intercept <- res2$intercept
165 slope <- res2$slope

```

166 **Analysis models**

167 Similarly, `pcr_analyze` provides two methods to model the C_T values and calculates the relative expression
 168 of target genes. 'delta_delta_ct' performs the $\Delta\Delta C_T$ method described previously. The average relative
 169 expression of the target gene in the condition of interest is given by the equations 1 & 2 and the
 170 standard deviation by 3. The calculations are applied to the dataset 'ct1', shown in Table 2 and Figure
 171 2A. 'relative_curve' performs the relative standard curve quantification, average relative expression/amount
 172 of the target gene in the condition of interest is given by equations 4 & 5 and the standard deviation by

Table 1. Average C_T value for c-myc and GAPDH at different input amounts

Input RNA (ng)	c-myc Average C_T	GAPDH Average C_T	ΔC_T c-myc - GAPDH
1.0	25.59 ± 0.04	22.64 ± 0.03	2.95 ± 0.05
0.5	26.77 ± 0.09	23.73 ± 0.05	3.04 ± 0.10
0.2	28.14 ± 0.05	25.12 ± 0.10	3.02 ± 0.11
0.1	29.18 ± 0.13	26.16 ± 0.02	3.01 ± 0.13
0.05	30.14 ± 0.03	27.17 ± 0.06	2.97 ± 0.07
0.02	31.44 ± 0.16	28.62 ± 0.10	2.82 ± 0.19
0.02	32.42 ± 0.12	29.45 ± 0.08	2.97 ± 0.14

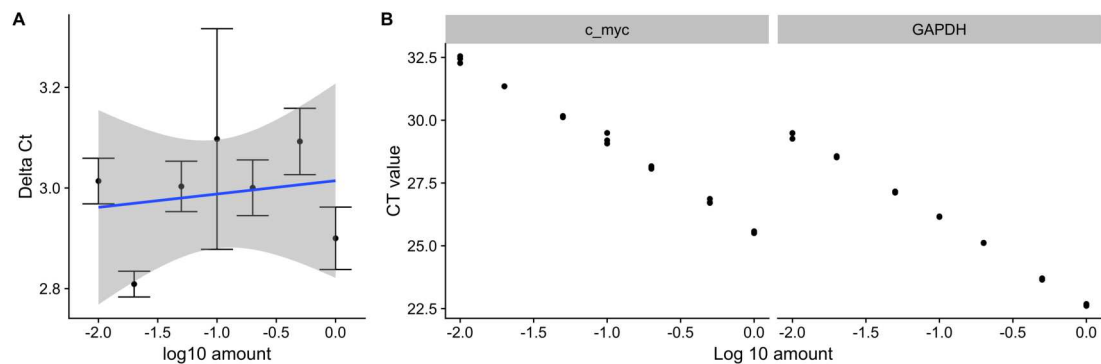


Figure 1. Amplification efficiency and standard curves of c-myc and GAPDH. Seven different dilutions of RNA were used as an input to synthesize cDNA, then to a real-time quantitative PCR reaction with c-myc and GAPDH primers. (A) ΔC_T values were calculated by subtracting the C_T values from three independent samples of the control gene(GAPDH) from the target c-myc. Averages and standard deviations are shown as points and error bars. The blue line represents the linear trend between the ΔC_T and log10 of the input amount. (B) C_T values from three independent samples of c-myc and GAPDH are shown with the corresponding log10 input amounts.

Table 2. Relative quantification using comparative ($\Delta\Delta C_T$) method (separate tubes)

Tissue	c-myc C_T	GAPDH C_T	ΔC_T c-myc - GAPDH	$\Delta\Delta C_T$ $\Delta C_T - \Delta C_{T,Brain}$	c-myc _N Rel. to Brain
Brain	30.72	23.7			
	30.34	23.56			
	30.58	23.47			
	30.34	23.65			
	30.5	23.69			
	30.43	23.68			
Average	30.49 ± 0.15	23.63 ± 0.09	6.86 ± 0.17	0.00 ± 0.17	1.0 (0.9–1.1)
Kidney	27.06	22.76			
	27.03	22.61			
	27.03	22.62			
	27.1	22.6			
	26.99	22.61			
	26.94	24.18			
Average	27.03 ± 0.06	22.66 ± 0.08	4.37 ± 0.10	-2.50 ± 0.10	5.6 (5.3–6.0)

173 equation 6 & 7. The calculation is applied to the same dataset 'ct1' and is shown in Table 3 and Figure
174 2B.

```

175 # pcr_analyze
176 ## locate and read raw ct data
177 fl <- system.file('extdata', 'ct1.csv', package = 'pcr')
178 ct1 <- read.csv(fl)
179
180 ## add grouping variable
181 group_var <- rep(c('brain', 'kidney'), each = 6)
182
183 # calculate all values and errors in one step
184 ## mode == 'separate_tube' default
185 res1 <- pcr_analyze(ct1,
186                     group_var = group_var,
187                     reference_gene = 'GAPDH',
188                     reference_group = 'brain')
189
190 ## calculate standard amounts and error
191 res2 <- pcr_analyze(ct1,
192                     group_var = group_var,
193                     reference_gene = 'GAPDH',
194                     reference_group = 'brain',
195                     intercept = intercept,
196                     slope = slope,
197                     method = 'relative_curve')
```

198 **Significance Testing**

199 Finally, `pcr_test` can be used to calculate useful statistics, p-values and confidence intervals on the
200 previous models. Two tests are available of the two-group comparisons; 't.test' and 'wilcox.test' to test
201 for the difference of the normalized target gene expression (ΔC_T) in one condition to another. Linear
202 regression, 'lm', can be applied to estimate these differences between multiple conditions and a reference
203 (Equation 8). The following code applies different testing methods to the dataset 'ct4'. The dataset was
204 published original in (Yuan et al., 2006), along with results of different testing method (Table 4). Table 5
205 shows the results of the three different tests as implemented in `pcr_test`.

Table 3. Relative quantification using the standard curve method (separate tube)

Tissue	c-myc (ng)	GAPDH (ng)	c-myc _N norm. to GAPDH	c-myc _N Rel. to Brain
Brain	0.033	0.51		
	0.043	0.56		
	0.036	0.59		
	0.043	0.53		
	0.039	0.51		
	0.040	0.52		
Average	0.039 ± 0.004	0.54 ± 0.034	0.07 ± 0.008	1.0 ± 0.12
Kidney	0.40	0.96		
	0.41	1.06		
	0.41	1.05		
	0.39	1.07		
	0.42	1.06		
	0.43	0.96		
Average	0.41 ± 0.016	1.02 ± 0.052	0.40 ± 0.025	5.5 ± 0.35

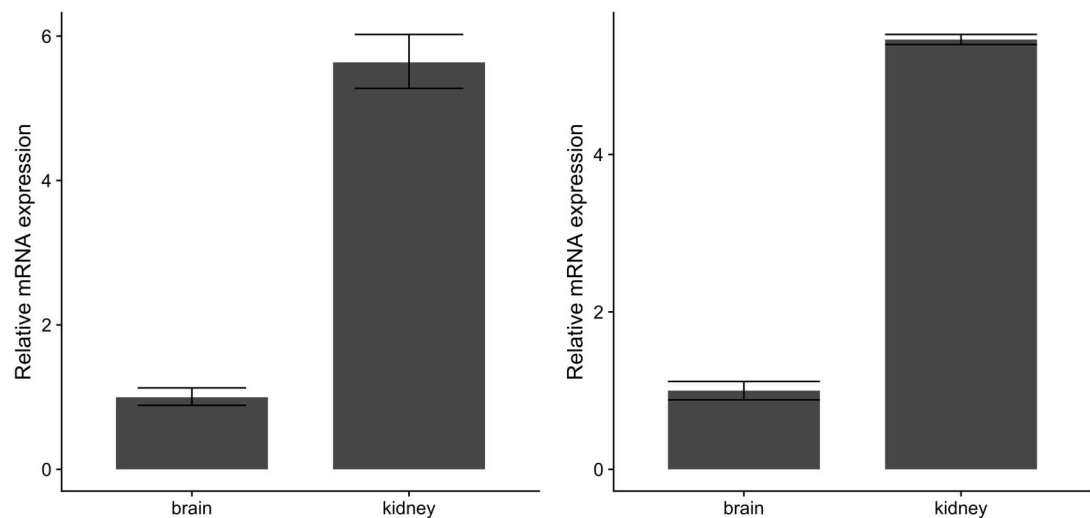


Figure 2. Relative expression of c-myc in human brain and kidney tissues. Total RNA from human brain and kidney tissues were used to synthesize cDNA and real-time quantitative PCR reaction with c-myc and GAPDH primers. C_T values from six independent replicates were used to calculate the expression of c-myc in the kidney normalized by GAPDH and relative to the brain using The $\Delta\Delta C_T$ (A) and the standard curve methods (B). Averages and standard deviations are shown as bars and error bars.

Table 4. Statistical significance using different testing methods

Test	$\Delta\Delta C_T$ (estimate)	<i>p</i> -value	Confidence Interval
Multiple Regression	-0.6848	<0.0001	(-0.4435, -0.9262)
ANOVA	-0.6848	<0.0001	(-0.4435, -0.9262)
<i>t</i> -test	-0.6848	<0.0001	(-0.4147, -0.955)
Wilcoxon test	-0.6354	<0.0001	(-0.4227, -0.8805)

Table 5. Different testing methods applied to the same dataset.

	gene	estimate	p-value	lower	upper	term
t.test	target	-0.68	0.00	-0.96	-0.41	
wilcox.test	target	-0.64	0.00	-0.88	-0.42	
lm	target	-0.68	0.00	-0.95	-0.41	group_vartreatment

```

206 # pcr_test
207 # locate and read data
208 fl <- system.file('extdata', 'ct4.csv', package = 'pcr')
209 ct4 <- read_csv(fl)
210
211 # make group variable
212 group <- rep(c('control', 'treatment'), each = 12)
213
214 # test using t-test
215 tst1 <- pcr_test(ct4,
216                 group_var = group,
217                 reference_gene = 'ref',
218                 reference_group = 'control',
219                 test = 't.test')
220
221 # test using wilcox.test
222 tst2 <- pcr_test(ct4,
223                 group_var = group,
224                 reference_gene = 'ref',
225                 reference_group = 'control',
226                 test = 'wilcox.test')
227
228 # testing using lm
229 tst3 <- pcr_test(ct4,
230                 group_var = group,
231                 reference_gene = 'ref',
232                 reference_group = 'control',
233                 test = 'lm')

```

234 Comparison with other packages

235 Pabinger et al. (2014) surveyed the tools used to analyze qPCR data across different platforms. They
 236 included 9 R packages which provide very useful analysis and visualization methods. Some of these
 237 packages focuses one certain model and some are designed to handle high-throughput qPCR data. Most
 238 of these packages are hosted in CRAN and a few on the Bioconductor so they adhere to Bioconductor
 239 methods and data containers. In comparison, pcr provides a unified interface for different quality
 240 assessment, analysis and testing models. The input and the output are tidy **data.frame**, and the package
 241 source code follows the tidyverse practices. This package targets the small scale qPCR experimental data

242 and the R user practitioners. The interface and documentation choices were made with such users in mind
243 and require no deep knowledge in specific data structures or complex statistical models.

244 **Limitations & future directions**

245 The current version of the pcr package (1.1.0) provides only methods to estimate the expression of genes
246 in a certain condition relative to another. Other methods were proposed for absolute quantification of the
247 copy number of genes in samples (Whelan et al., 2003). Also, the comparative C_T methods assume that
248 the PCR reaction has a close to perfect amplification rates. Other methods were proposed to model the
249 data when this assumption is not true (Liu and Saint, 2002; Tichopad et al., 2003). We are planning to
250 implement methods for absolute quantification and dealing with less than perfect amplification efficiency
251 cases in a future version of the package.

252 **CONCLUSION**

253 To sum, the pcr package is an open source R package for quality assessing, modeling and testing real-time
254 quantitative PCR data. The package provide an intuitive and unified interface for its main functions to
255 allow biologist to perform all necessary steps of qPCR analysis and produce graphs in a uniform way.

256 **ACKNOWLEDGMENTS**

257 We thank all lab members for the critical discussion at the development of this R package. This work was
258 supported by the Basic Science Research Program through the National Research Foundation of Korea
259 funded by the Ministry of Education (2015R1D1A01019753) and by the Ministry of Science, ICT and
260 Future Planning (NRF-2015R1A5A2008833).

261 **REFERENCES**

- 262 Bustin, S. A., Benes, V., Garson, J. A., Hellemans, J., Huggett, J., Kubista, M., Mueller, R., Nolan, T.,
263 Pfaffl, M. W., Shipley, G. L., Vandesompele, J., and Wittwer, C. T. (2009). The MIQE Guidelines:
264 Minimum Information for Publication of Quantitative Real-Time PCR Experiments. *Clinical Chemistry*,
265 55(4):611–622.
- 266 Bustin, S. A. and Nolan, T. (2004). Pitfalls of quantitative real-time reverse-transcription polymerase
267 chain reaction. *Journal of biomolecular techniques : JBT*, 15(3):155–66.
- 268 Liu, W. and Saint, D. A. (2002). A New Quantitative Method of Real Time Reverse Transcription
269 Polymerase Chain Reaction Assay Based on Simulation of Polymerase Chain Reaction Kinetics.
270 *Analytical Biochemistry*, 302(1):52–59.
- 271 Livak, K. J. and Schmittgen, T. D. (2001). Analysis of Relative Gene Expression Data Using Real-Time
272 Quantitative PCR and the Double Delta CT Method. *Methods*, 25(4):402–408.
- 273 Pabinger, S., Rödiger, S., Kriegner, A., Vierlinger, K., and Weinhäusel, A. (2014). A survey of tools for
274 the analysis of quantitative PCR (qPCR) data. *Biomolecular Detection and Quantification*.
- 275 Tichopad, A., Dilger, M., Schwarz, G., and Pfaffl, M. W. (2003). Standardized determination of real-time
276 PCR efficiency from a single reaction set-up. *Nucleic acids research*, 31(20):e122.
- 277 Whelan, J. A., Russell, N. B., and Whelan, M. A. (2003). A method for the absolute quantification of
278 cDNA using real-time PCR. *Journal of immunological methods*, 278(1-2):261–9.
- 279 Yuan, J., Reed, A., Chen, F., and Stewart, C. N. (2006). Statistical analysis of real-time PCR data. *BMC*
280 *Bioinformatics*, 7(1):85.