

*Improved performance of crystal structure solution from powder diffraction data through parameter tuning of a simulated annealing algorithm*

Article

Published Version

Kabova, E. A., Cole, J. C., Korb, O., López-Ibáñez, M., Williams, A. ORCID: <https://orcid.org/0000-0003-3654-7916> and Shankland, K. ORCID: <https://orcid.org/0000-0001-6566-0155> (2017) Improved performance of crystal structure solution from powder diffraction data through parameter tuning of a simulated annealing algorithm. *Journal of Applied Crystallography*, 50. pp. 1411-1420. ISSN 0021-8898 doi: 10.1107/S1600576717012602 Available at <https://centaur.reading.ac.uk/72695/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1107/S1600576717012602>

Publisher: Wiley

copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

[www.reading.ac.uk/centaur](http://www.reading.ac.uk/centaur)

## **CentAUR**

Central Archive at the University of Reading

Reading's research outputs online



## Improved performance of crystal structure solution from powder diffraction data through parameter tuning of a simulated annealing algorithm

**Elena A. Kabova, Jason C. Cole, Oliver Korb, Manuel López-Ibáñez, Adrian C. Williams and Kenneth Shankland**

*J. Appl. Cryst.* (2017). **50**, 1411–1420



**IUCr Journals**  
CRYSTALLOGRAPHY JOURNALS ONLINE

Copyright © International Union of Crystallography

Author(s) of this paper may load this reprint on their own web site or institutional repository provided that this cover page is retained. Republication of this article or its storage in electronic databases other than as specified above is not permitted without prior permission in writing from the IUCr.

For further information see <http://journals.iucr.org/services/authorrights.html>

# Improved performance of crystal structure solution from powder diffraction data through parameter tuning of a simulated annealing algorithm

Elena A. Kabova,<sup>a\*</sup> Jason C. Cole,<sup>b</sup> Oliver Korb,<sup>b</sup> Manuel López-Ibáñez,<sup>c</sup> Adrian C. Williams<sup>d</sup> and Kenneth Shankland<sup>d</sup>

Received 10 July 2017  
Accepted 1 September 2017

Edited by Th. Proffen, Oak Ridge National Laboratory, USA

**Keywords:** crystal structure determination; powder diffraction; simulated annealing; parameter tuning.

**Supporting information:** this article has supporting information at journals.iucr.org/j

<sup>a</sup>School of Pharmacy, University of Reading, Whiteknights Campus, Reading, Berks RG6 6AD, UK, <sup>b</sup>Cambridge Crystallographic Data Centre, 12 Union Road, Cambridge CB2 1EZ, UK, <sup>c</sup>Decision and Cognitive Sciences Research Centre, Alliance Manchester Business School, Booth Street East, Manchester, Lancashire M13 9SS, UK, and <sup>d</sup>School of Chemistry, Food and Pharmacy, University of Reading, Whiteknights, Reading, Berkshire RG6 6AP, UK.

\*Correspondence e-mail: e.kabova@reading.ac.uk

Significant gains in the performance of the simulated annealing algorithm in the *DASH* software package have been realized by using the *irace* automatic configuration tool to optimize the values of three key simulated annealing parameters. Specifically, the success rate in finding the global minimum in intensity  $\chi^2$  space is improved by up to an order of magnitude. The general applicability of these revised simulated annealing parameters is demonstrated using the crystal structure determinations of over 100 powder diffraction datasets.

## 1. Introduction

*DASH* (David *et al.*, 2006, 1998), a computer program for crystal structure determination from powder diffraction data (SDPD) which utilizes a simulated annealing (SA) algorithm, has previously been adapted to run on multiple CPU-core computers *via MDASH* (Griffin *et al.*, 2009b), distributed computing systems *via GDASH* (Griffin *et al.*, 2009a) and cloud computing systems *via CDASH* (Spillman *et al.*, 2015). Since its launch in 1999, the key SA control parameter values have remained unchanged; with good SDPD performance (Shankland *et al.*, 2013), there has been little incentive to vary them. However, as the program is applied to ever more complex structures, the chances of determining the correct crystal structure from any given SA run fall dramatically (Kabova, 2016). It is interesting and valuable to assess whether better parameterization of the simulated annealing algorithm can lead to better performance.

The effect upon *DASH* performance of individually varying the initial SA temperature ( $T_0$ ) and the cooling rate (CR) of the SA algorithm has previously been investigated (Shankland *et al.*, 2002), though only against a single powder X-ray diffraction dataset. The results showed that the automatic temperature setting in *DASH* was very effective but that setting the CR value too high (0.3, where the default value of CR is 0.02) halved the success rate<sup>1</sup> in locating the global minimum. The variation of the parameters  $N_1$  and  $N_2$ , which control the allocation of the available SA moves, was not investigated.

<sup>1</sup> Defined here as the number of SA runs that locate the global minimum, divided by the total number of SA runs performed, then expressed as a percentage.



Finding appropriate control parameter values is a challenge for all algorithm developers. Values may be adjusted manually, but such manual parameterization is easily subject to human bias, unless performed very carefully. For example, influenced by previously reported *DASH* results, it is highly unlikely that a researcher would choose to explore high values of CR during parameterization, because of (a) the expectation that it will lead to decreased success rates and (b) a desire to keep the number of SA parameter combinations being explored small, to save computational time. Automatic tuning algorithms ('tuners'), on the other hand, can implement the optimization using approaches which do not require the parameter space to be exhaustively explored and which alleviate the problems associated with human bias in parameter variation. The design and application of tuners is a dynamic area of research; examples include the work of Eiben & Smit (2012) on tuning evolutionary algorithms, use of SA for the optimization of mapping on network chips (Yang *et al.*, 2012), mixed integer programming (Hutter *et al.*, 2010) and general-purpose optimization algorithms (Balaprakash *et al.*, 2007).

Here, we report the use of the program *irace* (López-Ibáñez *et al.*, 2016), which implements the iterated racing procedure (Balaprakash *et al.*, 2007), to carry out SA parameter optimization of *DASH* against a training set of 40 powder diffraction datasets. *irace* has been shown (Pérez Cáceres *et al.*, 2014) to be well suited to tuning general-purpose algorithms with relatively large numbers of configurable parameters of different types, such as ordered, continuous, categorical and integer parameters, and has been applied to tune computer programs for optimization, machine learning and robotics (López-Ibáñez *et al.*, 2016).

Success rates in solving crystal structures from a library of over 100 powder diffraction datasets were then obtained for both the best performing SA parameter configuration suggested by *irace* and the default SA parameter values in *DASH*.

## 2. Materials and methods

### 2.1. Selection and composition of powder X-ray diffraction datasets

A study carried out by Florence *et al.* (2005) on 35 industrially relevant molecules concluded that crystal structures with greater than 20 degrees of freedom (DoF) could be classed as 'complex' and were broadly representative of the perceived limits of SDPD at the time. For the purposes of our work, a dataset size of 100 molecules was considered sufficiently large to facilitate an up-to-date, comprehensive and systematic study of the performance of *DASH*. The detailed selection criteria for dataset assembly are described fully elsewhere (Kabova, 2016), but the key criteria were as follows: (a) that the diffraction data should be derived from small organic molecules whose crystal structures have previously been solved from powder X-ray diffraction data, to ensure relevance and to permit assessment of the quality of the SA solutions; (b) that the crystal structures spanned a large DoF

**Table 1**  
Summary of the software used in this work.

Software	Version	Application	Reference
<i>DASH</i>	3.3.2	Indexing <sup>†</sup> Space group determination <sup>‡</sup> Pawley refinement Structure solution	David <i>et al.</i> (2006)
<i>dash.x</i>	3.3.2	Structure solution (under Linux)	CCDC (personal communications)
	3.3.1	<i>irace</i> calculations (under Linux)	
<i>MDASH</i>	3.1	Structure solution	Griffin <i>et al.</i> (2009b)
<i>TOPAS</i>	4.2	Indexing Pawley refinement Rietveld refinement	Coelho (2003)
CSD	5.36	Model building	Allen (2002)
<i>MarvinSketch</i>	6.0.5	Model building	ChemAxon (2011)
<i>ConQuest</i>	1.17	Structure mining of CSD	Bruno <i>et al.</i> (2002)
<i>Mercury</i>	3.3	Structure visualization	Macrae <i>et al.</i> (2008)
<i>Mogul</i>	1.6	Structure verification	Bruno <i>et al.</i> (2004)
<i>enCIFer</i>	1.51	CIF verification	Allen <i>et al.</i> (2004)
<i>PLATON</i>	1.51	Unit-cell conversion	Spek (2003)
<i>irace</i>	1.4	Algorithm optimization	López-Ibáñez <i>et al.</i> (2016)
<i>Minitab</i>	17.1.0.0	Statistical analysis	Minitab (2010)

<sup>†</sup> Via interface to *DICVOL91* (Boultif & Louër, 1991). <sup>‡</sup> With *ExtSym* as implemented in *DASH*.

range and that there should be multiple representatives of most of the DoF values.

In total, diffraction data associated with 101 crystal structures were assembled. To satisfy the requirements of *irace*, these were divided into two subsets – the 'training' set (A1–A40) and the 'test' set (B1–B61). The training set was a representative sample of 40 structures which, in the parameter-tuning experiments, was used by *irace* to optimize the SA parameters of *DASH*. The remaining 61 structures constituted the test set, which was then used to independently validate the performance of these optimized SA parameters.

The composition of the full dataset, in terms of complexity, can be summarized as follows: 50 structures with DoF < 14, 32 structures with 14 ≤ DoF ≤ 20, 16 structures with 21 ≤ DoF ≤ 30 and 3 structures with DoF > 30. Full molecular and crystallographic details for each structure are given in Table S1, Table S2 and Fig. S1 of the supporting information.

### 2.2. Software and hardware

The software employed in this work is summarized in Table 1, whilst the hardware used is summarized in Table 2. Note that *irace* is implemented as an R-package (R Core Team, 2011), which was obtained from CRAN (Hornik, 2015).

### 2.3. *irace* operation

A full description of the *irace* package lies outside the scope of this article; one can be found in the work of López-Ibáñez *et al.* (2016) and the description here is limited to definitions of a

Table 2

Summary of the hardware used in this work.

PC	CPU	RAM	Operating system
1	Intel Core 2 Quad Q9400 (2.66 GHz)	4 GB	Windows 7 Enterprise (64 bit)
2	2 × Intel Xeon E5520 (2.270 GHz)	32 GB	Windows Server 2008 R2 Datacenter (64 bit)
3	2 × Intel Xeon E5-2630 v2 (2.60 GHz)	16 GB	Windows 7 Professional (64 bit)
4	2 × Intel Xeon E5-2630 (2.30 GHz)	16 GB	Windows 7 Enterprise (64 bit) Ubuntu 13.04 (32 bit)

number of *irace*-related terms (Table 3) that are essential to understanding the optimization of the SA parameters, plus a description of the *irace* operation in that context.

The *irace* package implements the iterated racing procedure, which is a method for offline parameter tuning. Offline tuning consists of two clearly defined stages: tuning for identifying good performing configurations and testing (or deployment) of the best configurations found. The testing stage does not involve *irace*. Rather, it consists of evaluating the performance of each of the elite configurations suggested by *irace* against a set of instances (a test set) which were not included in the tuning stage.

**2.3.1. The *irace* procedure.** A single run of *irace* repeatedly iterates over three phases: (a) sampling new configurations (i.e. sets of SA parameters) according to a particular distribution, a truncated normal distribution in the case of numerical parameters; (b) selection of the best configurations by means of racing; and (c) updating the sampling distribution in order to bias future iterations towards optimal configurations. Racing is a well known method for the selection of the best candidate under uncertainty. In the context of algorithm configuration, candidate parameter configurations are evaluated over a sequence of training instances. As soon as there is enough evidence (for example, by means of the Friedman test) that some configurations are worse than the best one, the worst performing configurations are eliminated and the race continues until a minimum number of configurations remains or a maximum number of evaluations is reached. These three phases are repeated until some termination criterion is met – in the case of the current work, this is when the given budget (see Table 3 for definition) of *DASH* runs is reached. A representation of the workflow of SA parameter tuning is given in Fig. 1.

The number of iterations  $N_{\text{iterations}}$  performed during an *irace* run depends upon the number of optimizable parameters  $N_{\text{parameters}}$  and is calculated using

$$N_{\text{iterations}} = 2 + \log_2 N_{\text{parameters}}. \quad (1)$$

For optimization of the *DASH* SA parameters, each *irace* run comprises three iterations. Similarly, the budget for each iteration  $B_j$  is dependent on the total budget  $B$  and the number of iterations performed:

$$B_j = (B - B_{\text{used}})/(N_{\text{iterations}} - j + 1), \quad (2)$$

where  $j = 1, \dots, N_{\text{iterations}}$  and  $B_{\text{used}}$  is the sum of  $B_j$  for all previous iterations.

Table 3

Definitions of *irace*-related terms used in this work.

<i>irace</i> term	Symbol	Definition
Parameter space	$X$	The range of parameter values explored during the optimization
Tuning instance	$i$	A representative of the particular optimizable problem (e.g. crystal structures)
Training Set	n/a	A set of instances used in <i>irace</i> to benchmark the performance of <i>DASH</i>
Test Set	n/a	A set of instances unseen by <i>irace</i> , used to evaluate the <i>irace</i> results
Configuration	$\Theta_j$	A set of SA parameter values (e.g. CR = 0.02; $N_1 = 20$ ; $N_2 = 25$ )
Elite configuration	$\Theta_{\text{elite}}$	The best performing configuration, output at the end of an iteration
Experiment	n/a	An implementation of the algorithm with a specific configuration
Tuning budget	$B$	The maximum number of experiments (SA runs) performed
	$T^{\text{first}}$	The number of instances run before the first statistical test is applied
	$T^{\text{each}}$	The number of instances run before subsequent statistical tests are applied

Once the required inputs are in place, the first iteration (or ‘race’) starts with the uniform sampling of the parameter space and the generation of a set of parameter configurations,  $\Theta$ . For example:  $\Theta_1$  [CR = 0.20;  $N_1 = 6$ ;  $N_2 = 11$ ];  $\Theta_2$  [CR = 0.22;  $N_1 = 5$ ;  $N_2 = 25$ ];  $\dots$ ;  $\Theta_n$  (CR = 0.28;  $N_1 = 40$ ;  $N_2 = 31$ ).

Then the race is performed by following the steps given in Table 4. All subsequent iterations start with the generation of new candidate configurations based upon the elite configurations from the previous iteration. The number of candidate configurations generated at the start of an iteration reduces with the increasing number of iterations according to

$$\Theta_j = B_j/[\mu + \min(5, j)], \quad (3)$$

where  $\mu$  is a user-defined parameter (set to 5 in the current work), allowing control over the ratio between the budget and

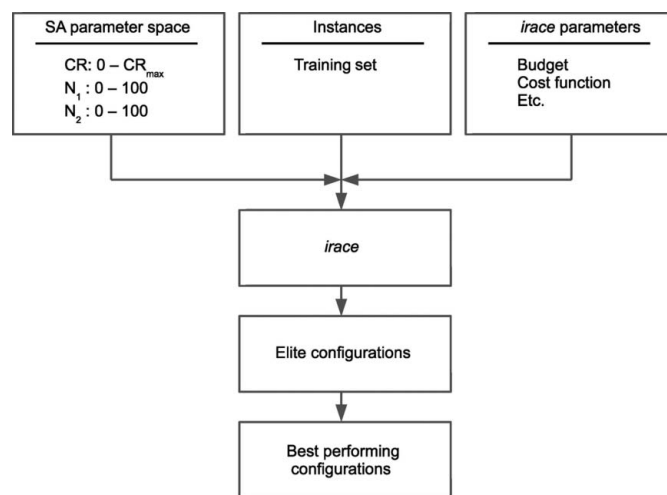


Figure 1

The SA parameter-tuning workflow. The *irace* ‘box’ represents the work carried out during the tuning stage. Once all cycles of the tuning are complete, the final elite configurations are output and carried over to the evaluation, which is performed independently of *irace*.



**Table 4**  
Steps performed during a race within a single *irace* iteration.

Step	<i>irace</i> step	Utilization of <i>DASH</i> by <i>irace</i>
1	Evaluate each candidate configuration on the first instance	Perform and evaluate <i>DASH</i> runs against instance 1, looping over the set of parameter configurations (configurations $j = 1 \dots n$ )
2	Continue the evaluation on subsequent instances until the number of instances reaches the predefined value of $T^{\text{first}}$	Using the same set of configurations ( $1 \dots n$ ), perform and evaluate <i>DASH</i> runs on these ( $T^{\text{first}}$ ) instances; $T^{\text{first}}$ was set to 5 in all experiments
3	Perform a statistical test on the evaluated configurations to identify statistically poorly performing configurations, if any	Check cost function values to determine which configurations resulted in the poorest <i>DASH</i> performance: for example configuration 1
4	Discard poorly performing configurations	Discard configuration 1
5	Run the next instance with the surviving configurations	Run the next instance with configurations $j = 2 \dots n$
6	Perform the statistical test every $T^{\text{each}}$ number of instances; predefined value	$T^{\text{each}}$ was set to 1, and thus statistical tests were performed after each instance
7	Continue until remaining budget is insufficient to test all remaining configurations on another instance ( $B_i < N_j^{\text{Surviving}}$ )	Continue until the number of remaining configurations ( $N_j^{\text{Surviving}}$ ) is larger than the remaining allowed number of <i>DASH</i> runs
8	Rank the surviving configurations based on their cost function value	Rank the surviving configurations based on their cost function value
9	Output $\Theta_{\text{elite}}$ (the three best ones from the surviving configurations)	At most, three sets of best performing SA parameters are output, e.g. $\Theta_6$ [CR = 0.16; $N_1 = 23$ ; $N_2 = 62$ ], $\Theta_{29}$ [CR = 0.15; $N_1 = 21$ ; $N_2 = 46$ ], $\Theta_2$ [CR = 0.22; $N_1 = 51$ ; $N_2 = 25$ ]

the number of configurations. When the total budget is exhausted, *irace* terminates. The configurations that survive the last race are then output in an analogous fashion to step 9 in Table 4. These are the SA parameter configurations, which are then evaluated against the test set.

**2.3.2. Cost function.** The cost function used by *irace* in the case of SA parameter optimization in *DASH* is defined as  $100\chi_{\text{profile}}^2/\chi_{\text{target}}^2$ , where  $\chi_{\text{profile}}^2$  is the familiar powder profile  $\chi^2$  value returned at the end of an SA run and  $\chi_{\text{target}}^2$  is the profile  $\chi^2$  value obtained for the correct crystal structure. To establish  $\chi^2$  for a given structure, a rigid-body Rietveld refinement of the previously deposited crystal structure was performed with *DASH*. The  $\chi^2$  value from this refinement was assumed to be the lowest achievable by the SA, and as such was set to be  $\chi_{\text{target}}^2$ .

**2.3.3. *irace* experiments.** A total of 14 *irace* runs were performed; ten of these were performed on the full training set (A1–A40) and were used to assess the validity of initial parameter value bounds. However, with various budgets in the range 5000–30 000 they consumed 799 days of CPU time. To focus subsequent computational effort in the area where improvement is most valuable (the more challenging structures with DoF  $\geq 14$ , where success rates are known to fall off significantly), the remaining four *irace* runs were performed on datasets A18–A40 only. These remaining *irace* runs had a total budget of 30 000 runs and took 516 days of CPU time. In all runs, the only optimizable parameters were CR,  $N_1$  and  $N_2$ . CR was varied as a real number in the range 0–0.3, whilst  $N_1$  and  $N_2$  were varied as integers in the range 0–100. These 14 *irace* runs alone utilized a total of 225 000 individual *DASH* runs and a total of 1315 CPU days.

## 2.4. Baseline *DASH* performance with default SA parameters

Initially, 50 SA runs were executed on all 101 structures, using the default *DASH* SA parameters (CR = 0.02;  $N_1 = 20$ ;

$N_2 = 25$ ). Each run was set to perform  $1 \times 10^7$  SA moves followed by a short simplex calculation. A  $\chi^2$  multiplier of 1 (CCDC, 2017) ensured the full number of SA moves was always carried out and that the SA was not terminated prematurely. The starting molecular conformers were randomly generated and all variable torsion angles were allowed to rotate freely (i.e. in the range 0–360°) during the SA calculations. Successful solutions were identified on the basis of their  $\chi^2$  value and further confirmed by comparison of coordinates with the reference crystal structure. The four crystal structures for which no reference structures had been previously deposited (A4, A6, B23 and B58) were considered solved when a favourable value of the  $\chi^2$  ratio (typically  $2 < \chi_{\text{profile}}^2 < 10$ , associated with a crystallographically sensible crystal structure) had been achieved ( $\chi_{\text{Ratio}}^2 = \chi_{\text{Profile}}^2/\chi_{\text{Pawley}}^2$ ,  $\chi_{\text{Pawley}}^2$  being the best  $\chi^2$  achieved by a Pawley-type fit to the data in question). A March–Dollase correction was introduced in the SA process for some structures (Table S2), in order to take account of intensity distortions attributable to preferred orientation of the crystallites in the samples.

For crystal structures that were not solved with the initial 50 SA runs, an additional 100 SA runs of  $1 \times 10^7$  SA moves were performed. If a structure remained unsolved after this further set of runs, a final attempt at a crystal structure solution was performed with another 500 SA runs of  $5 \times 10^7$  SA moves. In order to speed up these longer calculations, the 500 runs were performed using *MDASH* to spread the calculations over ten CPU cores. Those structures which still remained unsolved were considered to have a 0% success rate.

Whilst crystal structures deposited in the Cambridge Structural Database (CSD; Groom *et al.*, 2016) were used as the starting point for *Z*-matrix (Shankland, 2005) generation for the majority of the *DASH* calculations, the starting values of the flexible torsion angles were always randomized by *DASH* and so no advantage (other than the use of good quality bond lengths and bond angles) is conferred by this approach. Indeed, it represents the recommended approach in

Table 5

Elite configurations returned by *irace* calculations 11–14.

Run No.	Budget	Elite sets of SA parameters CR/ $N_1$ / $N_2$	CPU time (days)
11	30 000	0.27/59/50 0.25/31/56 0.28/63/51	93
12	30 000	0.25/75/29 0.27/70/25 0.26/74/23	115
13	30 000	0.25/46/62 0.30/35/69 0.29/38/57	199
14	30 000	0.24/18/84 0.21/16/85 0.21/19/91	109

global-optimization-based SDPD, of using the most accurate starting model that is available.

### 2.5. DASH performance with best performing configurations

The 12 elite configurations (Table 5) suggested by the four *irace* runs performed on complex structures with DoF  $\geq 14$  were initially performance tested against a representative subset of structures, consisting of A20, A25, A28, A29, A30, A32, A34, A38, B34, B44, B47, B48, B52 and B55, to manage the computational requirements (see footnote 2 for an explanation of why structures from the training set were included). Some additional, minor, manual variants on these configurations (Table 6) were also tested against these structures.

The six best performing SA parameter configurations (CR/ $N_1$ / $N_2$  = 0.27/73/56, 0.27/73/61, 0.27/73/51, 0.27/60/63, 0.25/35/86 and 0.25/46/62) from these tests were then tested against all molecules in the dataset.<sup>2</sup> The *DASH* runs performed using each configuration mirrored those of the *DASH* baseline calculations, *i.e.* initially 50 SA runs of  $1 \times 10^7$  moves were performed for all molecules, followed by 100 SA runs of  $1 \times 10^7$  moves for the unsuccessful examples. Finally, 500 SA runs of  $5 \times 10^7$  moves were carried out if required. To facilitate the direct comparison of results, all *DASH* runs were performed in an identical manner to those of the baseline, *i.e.* identical molecular models were used for the generation of Z-matrices, all variable torsion angles were allowed to rotate freely (*i.e.* in the range 0–360°) during the SA calculations, the same random seed values were used and a value of one was selected for the  $\chi^2$  multiplier to ensure that all SA moves were executed.

### 2.6. DASH performance analysis

The key performance indicator chosen is that of the success rate (SR), *i.e.* the percentage of any given set of SA runs that

<sup>2</sup> Whilst it is not standard practice to test the new configurations against the training set, it was done here to confirm that similarly improved performance was returned for structures that had been ‘seen’ by *irace* (training set) and ‘not seen’ by *irace* (test set).

Table 6

Additional configurations tested.

SA parameters CR/ $N_1$ / $N_2$			
0.27/73/61	0.27/59/63	0.25/35/86	0.19/20/73
0.27/73/56	0.27/53/61	0.25/31/86	0.19/73/20
0.27/73/51	0.27/53/51	0.25/31/76	0.19/25/63
0.27/73/41	0.27/49/40	0.25/31/66	0.19/63/25

successfully solve the crystal structure. Percentage values are then easily plotted against the number of DoF present in the structure, in order that general trends can be assessed. To facilitate comparison between baseline *DASH* performance and the performance of *DASH* using parameter configurations suggested by *irace*, an analysis based on the empirical log-of-the-odds (ELO) transform was performed. The ELO, as described by Cox & Snell (1989), takes the form given in equation (4):

$$\text{ELO} = \ln \left( \frac{r_i + 0.5}{n_i - r_i + 0.5} \right), \quad (4)$$

where  $i$  is the subject (*i.e.* each of the individual structures in the dataset),  $n_i$  is the maximum value of the sample (in this case the maximum SR, *i.e.* 100%) and  $r_i$  is the error associated with it, *i.e.* the actual SR value achieved. As such, equation (4) can be rewritten as

$$\text{ELO} = \ln \left( \frac{\text{SR}_i + 0.5}{100 - \text{SR}_i + 0.5} \right). \quad (5)$$

Regression analysis on the log-transformed data was performed using *Minitab* (Minitab, 2010).

## 3. Results

### 3.1. irace configurations

The elite configurations output by *irace* runs 11–14 are summarized in Table 5.

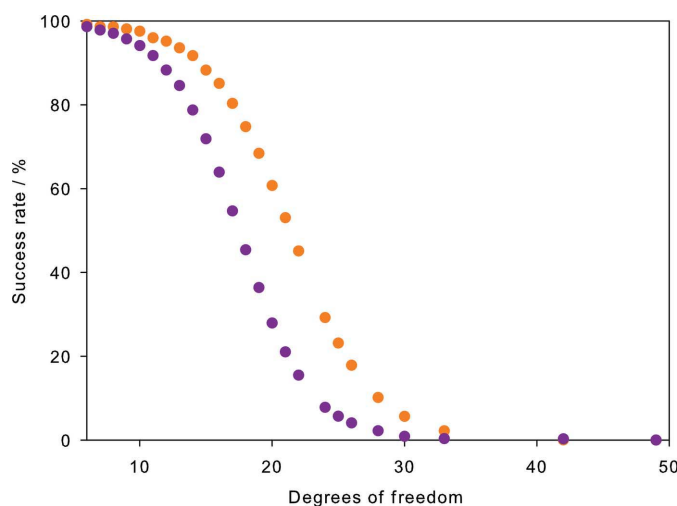


Figure 2

A comparison of the default (purple) and best performing (orange) SA parameter configuration models based on the ELO regressions.



Table 7

A comparison of SRs achieved using default and optimized SA parameter configurations.

No.	0.02/20/25	0.27/73/56	No.	0.02/20/25	0.27/73/56
A1	100	100	B12	96	100
A2	100	100	B13	100	100
A3	100	100	B14	100	100
A4	100	100	B15	66	98
A5	100	100	B16	100	100
A6	100	100	B17	100	100
A7	48	78	B18	70	100
A8	100	100	B19	100	100
A9	100	100	B20	100	100
A10	100	100	B21	44	60
A11	100	100	B22	100	100
A12	100	100	B23	98	100
A13	78	98	B24	96	98
A14	96	100	B25	100	100
A15	100	100	B26	84	98
A16	42	74	B27	44	78
A17	100	100	B28	100	100
A18	4	6	B29	92	100
A19	14	12	B30	64	98
A20	34	88	B31	58	50
A21	56	78	B32	100	100
A22	28	74	B33	100	100
A23	54	92	B34	50	100
A24	50	84	B35	14	48
A25	2†	24	B36	4	12
A26	1†	10	B37	12	30
A27	78	96	B38	36	76
A28	8	40	B39	4	14
A29	60	96	B40	8	26
A30	34	56	B41	98	100
A31	16	20	B42	20	44
A32	18	54	B43	12	32
A33	14	40	B44	8	48
A34	4	36	B45	14	54
A35	14	48	B46	4	70
A36	46	72	B47	14	54
A37	0‡	1†	B48	4†	12
A38	98	100	B49	0‡	1†
A39	1†	4	B50	0.2‡	0.4‡
A40	0.2‡	4	B51	1†	6
B1	92	100	B52	9.4‡	18
B2	100	100	B53	2	22
B3	100	100	B54	2‡	1†
B4	100	100	B55	4	90
B5	100	100	B56	0‡	0‡
B6	100	100	B57	0‡	0.4‡
B7	100	100	B58	78	100
B8	100	100	B59	0‡	0.2‡
B9	100	100	B60	0.4‡	1†
B10	100	100	B61	0‡	0‡
B11	100	100			

† The reported SR is achieved with 100 SA runs (each performing  $1 \times 10^7$  SA steps). ‡ The reported SR is achieved with 500 SA runs (each performing  $5 \times 10^7$  SA steps).

### 3.2. DASH performance

The performance of *DASH*, using default SA parameters and the best performing SA configuration from the experiments outlined in §2.5, is summarized in Tables 7 and 8. Full details of the baseline performance are given in Tables S3 and S5 of the supplementary information, whilst details of the six best performing SA configurations are given in Tables S4 and S5. The calculations for this element of the work required just

Table 8

Average success rates for each DoF, for both the default and the best performing SA parameter sets.

DoF	Number of representatives in dataset	Average SR (%)	
		0.02/20/25	0.27/73/56
6	4	98.0	100
7	5	100	100
8	5	100	100
9	6	93.7	99.7
10	8	82.8	92.3
11	7	99.1	99.7
12	6	84.3	95.7
13	9	84.6	93.8
14	11	24.9	42.6
15	2	45.0	83.0
16	7	24.3	64.9
17	3	27.0	43.3
18	5	26.4	53.2
20	4	9.0†	26.8
21	2	0.6†	3.2†
22	1	4.0	36.0
24	3	23.1	46.0
25	2	1†	11.5
26	1	2†	1
28	5	20.6†	38.9†
30	2	39.1†	52.0
33	1	0.0†	0.2†
42	1	0.4†	1
49	1	0.0†	0.0†

† Value given includes success rates based on experiments that required 500 SA runs.

over 3348 days of CPU time. The ELO analysis of the baseline *DASH* performance yields

$$\text{ELO} = 6.565 - 0.375 \text{ DoF}_{\text{total}}, \quad (6)$$

with an  $R^2$  (where  $R^2$  = explained variation/total variation) of 53.73% and a  $p$  value of 0.00 for the total DoF,  $\text{DoF}_{\text{total}}$ , showing them to be a statistically significant factor in determining success rate. The ELO analysis of *DASH* performance using the best performing SA parameter set yields

$$\text{ELO} = 7.013 - 0.329 \text{ DoF}_{\text{total}}, \quad (7)$$

with an  $R^2$  of 51.7% and a  $p$  value of 0.00 for the  $\text{DoF}_{\text{total}}$ . Using the above equations, a predicted SR can be calculated, for any structure, based on the total DoF. For the current dataset of 101 crystal structures, the calculated *DASH* performance for both the default and the best performing SA configurations is shown in Fig. 2. The fit of each ELO model to the experimental data is shown in Fig. S2 of the supplementary information.

### 4. Discussion

The objective of this work was to determine if the performance of the simulated annealing algorithm in *DASH* could be improved through the optimization of three key SA parameters using the *irace* program. The results indicate that there was considerable room for improvement in performance over that obtained using the default SA parameters which have been in place since the initial release of *DASH*. This is most

Table 9

Distribution of space groups (transformed to standard settings) within the full dataset and the CSD.

Space group	No. of structures in dataset	No. of organic powder structures in CSD	% of structures in dataset	% of organic powder structures in CSD
$P2_1/c$	40	355	39.6	35.8
$P1$	23	185	22.8	18.6
$P2_12_12_1$	16	146	15.8	14.7
$P2_1$	10	100	9.9	10.1
$Pbca$	3	63	3.0	6.4
$Pna2_1$	2	33	2.0	3.3
$C2/c$	2	41	2.0	4.1
$P1$	1	22	1.0	2.2
$Pbc2_1$	1	20	1.0	2.0
$I2$	1	13	1.0	1.3
$Pbcn$	1	7	1.0	0.7
$Cmca$	1	7	1.0	0.7

clearly indicated in Fig. 2, which compares the best fit lines obtained by ELO analysis of baseline performance and the performance of the best performing SA configuration returned as a result of the *irace* experiments. The marked shift to the right seen for the best performing configuration demonstrates the significant gains in success rate and consequent ability to tackle more complex structures in a finite time period. This improvement, its determination, its range of applicability and its significance are discussed more fully below.

#### 4.1. Structural complexity and the dataset composition

A recent analysis of structural complexity of crystal structures in the CSD (Shankland *et al.*, 2013) showed that the average complexity of deposited structures since the year 2000 is approximately 52 atoms in the asymmetric unit and approximately 13 DoF, and showed that SDPD methods are well placed to address problems of such complexity. The analysis also showed SDPD to be capable of solving problems of much greater than average complexity.

The structures in the dataset employed in this current work span a wide complexity range ( $6 \leq \text{DoF} \leq 49$ , with 50% of structures having  $\text{DoF} \geq 14$ ) and were chosen to ensure that any improved performance is directly relevant to structures that are likely to be attempted by SDPD now or in the near future. Such applications include the following: single molecules, salts, hydrates, solvates and organometallic structures; rigid molecules and conformationally flexible molecules ( $0 \leq \text{DoF}_{\text{torsional}} \leq 43$ ); cases with  $0.5 \leq Z' \leq 4$ ; laboratory-based and synchrotron-based X-ray data; representative coverage of typically encountered space groups (see Table 9).

The resolution (minimum  $d$  spacing) of the powder data and the number of reflections used in the Pawley refinement are two fundamental factors expected to influence both the SR and the quality of the *DASH* solution. Large variations of those factors were observed within the dataset, with B3 having the lowest resolution (only 3.64 Å) and only 19 contributing reflections.

#### 4.2. Considerations in setting up the *irace* runs

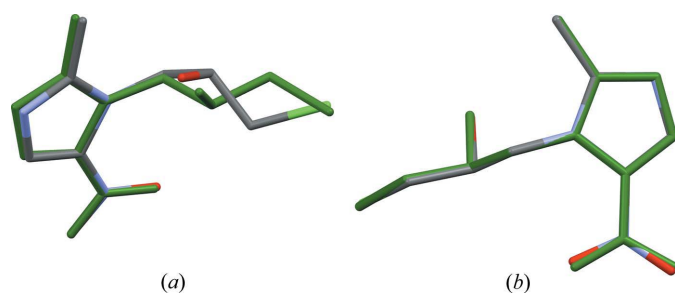
*DASH* has few user-controllable parameters which affect its performance: the starting temperature ( $T_0$ ), the cooling rate (CR), and the integers  $N_1$  and  $N_2$ , whose product ( $N_1N_2$ ) governs the number of SA moves performed at each temperature before a cooling step is applied.<sup>3</sup> Currently a value of '0' is the default for  $T_0$ , which instructs *DASH* to automatically determine an optimal value of this parameter for the structure under investigation. This is achieved by performing a short preliminary SA run during which the variation in  $\chi^2$  at different temperatures is examined. The temperature above which no significant variations in the  $\chi^2$  values are observed is selected as the appropriate starting temperature for the SA. The ranges explored for the remaining parameters CR,  $N_1$  and  $N_2$  were set pragmatically, recognizing that they needed to accommodate significant changes from the default *DASH* parameter values but also acknowledging the computational demands of spanning large ranges of parameter space: large ranges require a large *irace* budget to ensure good coverage. There was some uncertainty as to whether *irace* runs of large numbers of *DASH* calculations (*i.e.* a larger budget) would give superior results to their small-budget counterparts. Runs of *irace* with large budgets (*e.g.* 30 000 *DASH* runs) were generally expected to give better results, owing to the larger number of evaluations carried out. Ultimately, to take account of the stochastic nature of *irace*, and to explore all options, *irace* runs of varying budgets were performed.

#### 4.3. Baseline versus optimized *DASH* performance

It is clear from Tables 7 and 8 that improvements in the SR are seen right across the dataset. Of particular interest are all compounds which, during the baseline calculations, required 500 SA runs to solve, *i.e.* A40, B50, B52, B54 and B60. All of these returned a solution within the first 100 SA runs with the best performing SA configuration, a remarkable improvement in performance. Furthermore, four of the six compounds that gave 0% SR with the default settings (A37, B49, B57 and B59) now returned a solution with the best performing configuration, with only structures B56 and B61 remaining unsolved.

Close examination of all results shows that performance gains are achieved not only by reducing the overall number of SA runs needed to reach a solution but also by reducing the number of SA moves required to reach a solution, especially for the 'complex' examples. Furthermore, the quality of the solutions obtained using the best performing SA configuration was always at least as good as that obtained using default settings and, for some of the complex structures, substantially better. For example, the root-mean-square deviation (RMSD) for a 15 molecule overlay in *Mercury* (Macrae *et al.*, 2008) for

<sup>3</sup>  $N_1$  and  $N_2$  are parameters which control the number of SA steps performed during the individual SA runs.  $N_1$  is the number of times each of the DoF is adjusted before the DoF step lengths are altered.  $N_2$  is the number of times that this cycle is repeated before a temperature reduction is applied. In other words,  $N_1N_2\text{DoF} = \text{total SA moves at a given temperature}$ .



**Figure 3**

Crystal structure overlay of the reference crystal structure for A40 (dark green) and (a) the best structure obtained using *DASH* with default settings (0.02/20/25; RMSD = 0.296 Å) and (b) the best structure obtained using *DASH* with the best performing SA parameter configuration (0.27/73/56; RMSD = 0.075 Å). For clarity, H atoms have been omitted and only one ornidazole molecule, which is representative of the goodness of fit for all three molecules in the asymmetric unit cell, is shown.

structure B60 fell from 0.498 to 0.225 Å and that of A40 fell from 0.296 to 0.075 Å (Fig. 3).

Table 7 emphasizes the wide variation in changes of success rate as a result of changing from the defaults to the best performing SA parameter set. For many relatively simple structures that can be solved 100% of the time with the default settings, there is clearly no room for improvement. For several very complex structures that cannot be solved with the defaults, there is clearly infinite room for improvement. Our estimate of an overall ‘order of magnitude’ improvement in performance is derived from datasets with more than 14 DoF, and by arbitrarily assigning an improvement factor of 100 for cases where the default success rate was zero whilst the best performing SA parameter set success rate was nonzero.

#### 4.4. ELO analysis

An ELO analysis has been chosen as suitable for modelling the S-shaped curves that describe the changes of SR as a function of DoF. This is most clearly indicated in Fig. 2, where the marked shift to the right seen for the best performing configuration demonstrates the significant gains in success rate and consequent ability to tackle more complex structures in a finite time period. The  $R^2$  values for the ELO fits to the data are not high, indicating that structure complexity is only one factor in determining the SR.

A more detailed analysis of the data shows that a better fit to the data is achieved when the individual components of the DoF (*i.e.* positional, orientational and torsional) are considered separately in the ELO analysis (Kabova, 2016). The resultant model better accounts for the high SR observed for some compounds with large numbers of positional DoF and highlights that, in general, structures with large numbers of torsional DoF are more difficult to solve than structures with the same total DoF, but higher numbers of non-torsional DoF.

#### 5. Conclusions and recommendations

The significance of the results presented here lies in the fact that a remarkable improvement in performance has been

**Table 10**

Recommendations for setting up *DASH* runs, based on problem complexity.

Complexity	No. of SA runs	No. of SA moves	Note
DoF < 14	50	$5 \times 10^6$	–
$14 \leq \text{DoF} \leq 20$	50	$1 \times 10^7$	–
$21 \leq \text{DoF} \leq 27$	100	$1 \times 10^7$	–
DoF > 27	500	$5 \times 10^7$	†

† The use of prior conformational knowledge (*e.g.* obtained from the Cambridge Structural Database via *Mogul*) is considered to be highly beneficial for structures of this complexity (Kabova, 2016).

achieved merely by adjusting the SA control parameter values, with no changes to the underlying SA algorithm. The contribution of *irace* (López-Ibáñez *et al.*, 2016) in deriving the best performing SA parameter configuration cannot be underestimated. It is unlikely that a set of control parameters which included such a high cooling rate would have been considered by a process of manual selection. Importantly, the best performing SA configuration can be utilized immediately by manually entering the appropriate parameter values into *DASH*, in the ‘SA options’ window. It seems probable that the approach taken in this work can be applied to many other crystallographic programs that rely upon optimization algorithms that have not themselves been optimized in terms of their performance with respect to key control parameters.

Finally, some general recommendations, based on the number of DoF in a crystal structure under investigation, can be made regarding the number of runs and SA moves required by *DASH* to give a high level of certainty that the crystal structure will be solved; these are listed in Table 10.

#### 6. Availability and documentation

Details of *DASH*'s availability can be found at <https://www.ccdc.cam.ac.uk/solutions/csd-materials/components/dash/>.

#### 7. Related literature

Details of the 101 crystal structures used in this work are reported in the supporting information. The related references are as follows: Albov *et al.* (2006), Alleaume (1967), Assaad & Rukiah (2011), Ávila *et al.* (2009), Bamgboye & Sowerby (1986), Bauer *et al.* (2001), Beale & Stephenson (1972), Bekö *et al.* (2012), Borea *et al.* (1987), Bortolotti *et al.* (2011), Brammer & Stevens (1989), Brüning *et al.* (2010), Burley (2005), Burley *et al.* (2006), Bushmarinov *et al.* (2012), Carpy *et al.* (1985), Chernyshev *et al.* (2000, 2002, 2010), Clegg & Teat (2000), David *et al.* (1998), Dinnebier *et al.* (2000), Donaldson *et al.* (1981), Dorokhov *et al.* (2007), Dupont & Dideberg (1972), Eibl *et al.* (2009), Fernandes *et al.* (2006), Fernandes, Florence *et al.* (2007a,b), Fernandes, Shankland *et al.* (2007), Florence *et al.* (2003, 2005, 2008), Freer *et al.* (1993), Fries *et al.* (1971), Fujinaga & James (1980), Gadret *et al.* (1976), Haynes *et al.* (2006), Helmholdt *et al.* (2002), Himes *et al.* (1981), Hodgson & Asplund (1991), Hulme *et al.* (2006), Ivashevskaja

*et al.* (2003), Ivashevskaya *et al.* (2009), Johnston *et al.* (2004), Kato *et al.* (1979), Kennedy *et al.* (2001), Kojicprodic *et al.* (1984), Koo *et al.* (1980), Lefebvre *et al.* (2005), Llinàs *et al.* (2006), Maccaroni *et al.* (2010), Majumder *et al.* (2013), Marder (2004), Nichols & Frampton (1998), Nishibori *et al.* (2008), Noguchi, Fujiki *et al.* (2012), Noguchi, Miura *et al.* (2012), Nowell *et al.* (2002), Post & Horn (1977), Rohlíček *et al.* (2010), Rukiah & Al-Ktaifani (2011), Rukiah & Assaad (2010), Rukiah *et al.* (2004), Schmidt *et al.* (2005), Sergeev *et al.* (2010), Shankland (personal communication), Shankland *et al.* (1996, 2001), Shanmuga Sundara Raj *et al.* (2000), Shin *et al.* (1995), Smrčok *et al.* (2007), Sorrenti *et al.* (2013), Steiner (2000), van de Streek *et al.* (2009), Vallcorba *et al.* (2011), Yatsenko *et al.* (2001).

## Acknowledgements

EAK thanks the University of Reading and the Cambridge Crystallographic Data Centre (CCDC) for funding. We thank Mark Spillman and David Edgeley for their help with various computational matters pertaining to the rapid execution of *DASH*, particularly under the Linux operating system. We are also grateful to the University of Reading Chemical Analysis Facility for powder X-ray diffraction facilities.

## References

- Albov, D. V., Jassem, A. & Kuznetsov, A. I. (2006). *Acta Cryst.* **E62**, o1449–o1451.
- Allen, F. H. (2002). *Acta Cryst.* **B58**, 380–388.
- Allen, F. H., Johnson, O., Shields, G. P., Smith, B. R. & Towler, M. (2004). *J. Appl. Cryst.* **37**, 335–338.
- Alleaume, M. (1967). PhD thesis, University of Bordeaux, France.
- Assaad, T. & Rukiah, M. (2011). *Acta Cryst.* **C67**, o469–o472.
- Ávila, E. E., Mora, A. J., Delgado, G. E., Contreras, R. R., Rincón, L., Fitch, A. N. & Brunelli, M. (2009). *Acta Cryst.* **B65**, 639–646.
- Balaprakash, P., Birattari, M. & Stützle, T. (2007). *Hybrid Metaheuristics*, Lecture Notes in Computer Science, Vol. 4771, edited by T. Bartz-Beielstein, M. Blesa Aguilera, C. Blum, B. Naujoks, A. Roli, G. Rudolph & M. Sampels, pp. 108–122. Berlin, Heidelberg: Springer.
- Bamgboye, T. T. & Sowerby, D. B. (1986). *Polyhedron*, **5**, 1487–1488.
- Bauer, J., Spanton, S., Henry, R., Quick, J., Dziki, W., Porter, W. & Morris, J. (2001). *Pharm. Res.* **18**, 859–866.
- Beale, J. P. & Stephenson, N. C. (1972). *J. Pharm. Pharmacol.* **24**, 277–280.
- Bekö, S. L., Urmann, D., Lakatos, A., Glaubitz, C. & Schmidt, M. U. (2012). *Acta Cryst.* **C68**, o144–o148.
- Borea, P. A., Gilli, G., Bertolasi, V. & Ferretti, V. (1987). *Mol. Pharmacol.* **31**, 334–344.
- Bortolotti, M., Lonardelli, I. & Pepponi, G. (2011). *Acta Cryst.* **B67**, 357–364.
- Boultif, A. & Louër, D. (1991). *J. Appl. Cryst.* **24**, 987–993.
- Brammer, L. & Stevens, E. D. (1989). *Acta Cryst.* **C45**, 400–403.
- Brüning, J., Alig, E. & Schmidt, M. U. (2010). *Acta Cryst.* **C66**, o341–o344.
- Bruno, I. J., Cole, J. C., Edgington, P. R., Kessler, M., Macrae, C. F., McCabe, P., Pearson, J. & Taylor, R. (2002). *Acta Cryst.* **B58**, 389–397.
- Bruno, I. J., Cole, J. C., Kessler, M., Luo, J., Motherwell, W. D. S., Purkis, L. H., Smith, B. R., Taylor, R., Cooper, R. I., Harris, S. E. & Orpen, A. G. (2004). *J. Chem. Inf. Comput. Sci.* **44**, 2133–2144.
- Burley, J. C. (2005). *Acta Cryst.* **B61**, 710–716.
- Burley, J. C., van de Streek, J. & Stephens, P. W. (2006). *Acta Cryst.* **E62**, o797–o799.
- Bushmarinov, I. S., Dmitrienko, A. O., Korlyukov, A. A. & Antipin, M. Yu. (2012). *J. Appl. Cryst.* **45**, 1187–1197.
- Carpy, A., Léger, J.-M. & Melchiorre, C. (1985). *Acta Cryst.* **C41**, 624–627.
- CCDC (2017). *DASH User Guide and Tutorials*, <https://www.ccdc.cam.ac.uk/support-and-resources/ccdcresources/dash.pdf>.
- ChemAxon (2011). *Marvin 5.4.1.1*, <http://www.chemaxon.com>.
- Chernyshev, V. V., Kukushkin, S. Y. & Velikodny, Y. A. (2010). *Acta Cryst.* **E66**, o613.
- Chernyshev, V. V., Tafeenko, V. A., Makarov, V. A., Sonneveld, E. J. & Schenk, H. (2000). *Acta Cryst.* **C56**, 1159–1160.
- Chernyshev, V. V., Yatsenko, A. V., Kuvshinov, A. M. & Shevelev, S. A. (2002). *J. Appl. Cryst.* **35**, 669–673.
- Clegg, W. & Teat, S. J. (2000). *Acta Cryst.* **C56**, 1343–1345.
- Coelho, A. (2003). *TOPAS User Manual*. Bruker AXS GmbH, Karlsruhe, Germany.
- Cox, D. R. & Snell, E. J. (1989). *The Analysis of Binary Data*, 2nd ed. London: Chapman and Hall.
- David, W. I. F., Shankland, K. & Shankland, N. (1998). *Chem. Commun.* pp. 931–932.
- David, W. I. F., Shankland, K., van de Streek, J., Pidcock, E., Motherwell, W. D. S. & Cole, J. C. (2006). *J. Appl. Cryst.* **39**, 910–915.
- Dinnebier, R. E., Sieger, P., Nar, H., Shankland, K. & David, W. I. F. (2000). *J. Pharm. Sci.* **89**, 1465–1479.
- Donaldson, J. D., Leary, J. R., Ross, S. D., Thomas, M. J. K. & Smith, C. H. (1981). *Acta Cryst.* **B37**, 2245–2248.
- Dorokhov, A. V., Chernyshov, D. Y., Burlov, A. S., Garnovskii, A. D., Ivanova, I. S., Pyatova, E. N., Tsivadze, A. Y., Aslanov, L. A. & Chernyshev, V. V. (2007). *Acta Cryst.* **B63**, 402–410.
- Dupont, L. & Dideberg, O. (1972). *Acta Cryst.* **B28**, 2340–2347.
- Eibl, S., Fitch, A., Brunelli, M., Evans, A. D., Pattison, P., Plazanet, M., Johnson, M. R., Alba-Simionesco, C. & Schober, H. (2009). *Acta Cryst.* **C65**, o278–o280.
- Eiben, A. E. & Smit, S. K. (2012). *Autonomous Search*, edited by Y. Hamadi, E. Monfroy & F. Saubion, pp. 15–36. Berlin, Heidelberg: Springer.
- Fernandes, P., Florence, A., Shankland, K., Karamertzanis, P. G., Hulme, A. T. & Anandamanoharan, P. (2007a). *Acta Cryst.* **E63**, o247–o249.
- Fernandes, P., Florence, A. J., Shankland, K., Karamertzanis, P. G., Hulme, A. T. & Anandamanoharan, R. P. (2007b). *Acta Cryst.* **E63**, o202–o204.
- Fernandes, P., Florence, A. J., Shankland, K., Shankland, N. & Johnston, A. (2006). *Acta Cryst.* **E62**, o2216–o2218.
- Fernandes, P., Shankland, K., Florence, A. J., Shankland, N. & Johnston, A. (2007). *J. Pharm. Sci.* **96**, 1192–1202.
- Florence, A. J., Baumgartner, B., Weston, C., Shankland, N., Kennedy, A. R., Shankland, K. & David, W. I. F. (2003). *J. Pharm. Sci.* **92**, 1930–1938.
- Florence, A. J., Shankland, K., Gelbrich, T., Hursthouse, M. B., Shankland, N., Johnston, A., Fernandes, P. & Leech, C. K. (2008). *CrystEngComm*, **10**, 26–28.
- Florence, A. J., Shankland, N., Shankland, K., David, W. I. F., Pidcock, E., Xu, X., Johnston, A., Kennedy, A. R., Cox, P. J., Evans, J. S. O., Steele, G., Cosgrove, S. D. & Frampton, C. S. (2005). *J. Appl. Cryst.* **38**, 249–259.
- Freer, A. A., Bunyan, J. M., Shankland, N. & Sheen, D. B. (1993). *Acta Cryst.* **C49**, 1378–1380.
- Fries, D. C., Rao, S. T. & Sundaralingam, M. (1971). *Acta Cryst.* **B27**, 994–1005.
- Fujinaga, M. & James, M. N. G. (1980). *Acta Cryst.* **B36**, 3196–3199.
- Gadret, M., Goursolle, M., Leger, J. M., Colleter, J. C. & Carpy, A. (1976). *Acta Cryst.* **B32**, 2757–2761.
- Griffin, T. A. N., Shankland, K., van de Streek, J. & Cole, J. (2009a). *J. Appl. Cryst.* **42**, 356–359.



- Griffin, T. A. N., Shankland, K., van de Streek, J. & Cole, J. (2009*b*). *J. Appl. Cryst.* **42**, 360–361.
- Groom, C. R., Bruno, I. J., Lightfoot, M. P. & Ward, S. C. (2016). *Acta Cryst. B* **72**, 171–179.
- Haynes, D. A., Van de Streek, J., Burley, J. C., Jones, W. & Motherwell, W. D. S. (2006). *Acta Cryst. E* **62**, o1170–o1172.
- Helmholdt, R. B., Peschar, R. & Schenk, H. (2002). *Acta Cryst. B* **58**, 134–139.
- Himes, V. L., Mighell, A. D. & De Camp, W. H. (1981). *Acta Cryst. B* **37**, 2242–2245.
- Hodgson, D. J. & Asplund, R. O. (1991). *Acta Cryst. C* **47**, 1986–1987.
- Hornik, K. (2015). *R FAQ*, <https://CRAN.R-project.org/doc/FAQ/R-FAQ.html>.
- Hulme, A. T., Fernandes, P., Florence, A., Johnston, A. & Shankland, K. (2006). *Acta Cryst. E* **62**, o3046–o3048.
- Hutter, F., Hoos, H. & Leyton-Brown, K. (2010). *Ann. Math. Artif. Intell.* **60**, 65–89.
- Ivashevskaja, S. N., Aleshina, L. A., Andreev, V. P., Nizhnik, Y. P., Chernyshev, V. V. & Schenk, H. (2003). *Acta Cryst. E* **59**, o1006–o1008.
- Ivashevskaya, S. N., van de Streek, J., Djanhan, J. E., Brüning, J., Alig, E., Bolte, M., Schmidt, M. U., Blaschka, P., Höffken, H. W. & Erk, P. (2009). *Acta Cryst. B* **65**, 212–222.
- Johnston, A., Florence, A. J., Shankland, K., Markvardsen, A., Shankland, N., Steele, G. & Cosgrove, S. D. (2004). *Acta Cryst. E* **60**, o1751–o1753.
- Kabova, E. A. (2016). PhD thesis, University of Reading, UK.
- Kato, Y., Haimoto, Y. & Sakurai, K. (1979). *Bull. Chem. Soc. Jpn.* **52**, 233–234.
- Kennedy, A. R., Hughes, M. P., Monaghan, M. L., Staunton, E., Teat, S. J. & Smith, W. E. (2001). *J. Chem. Soc. Dalton Trans.* pp. 2199–2205.
- Kojicprodic, B., Ruzictoros, Z., Sunjic, V., Decorte, E. & Moimas, F. (1984). *Helv. Chim. Acta*, **67**, 916–926.
- Koo, C. H., Cho, S. I. & Yeon, Y. H. (1980). *Arch. Pharm. Res.* **3**, 37–49.
- Lefebvre, J., Willart, J.-F., Caron, V., Lefort, R., Affouard, F. & Danède, F. (2005). *Acta Cryst. B* **61**, 455–463.
- Llinàs, A., Fábian, L., Burley, J. C., van de Streek, J. & Goodman, J. M. (2006). *Acta Cryst. E* **62**, o4196–o4199.
- López-Ibáñez, M., Dubois-Lacoste, J., Pérez Cáceres, L., Birattari, M. & Stützel, T. (2016). *Oper. Res. Perspect.* **3**, 43–58.
- Maccaroni, E., Malpezzi, L. & Masciocchi, N. (2010). *Acta Cryst. E* **66**, o2511.
- Macrae, C. F., Bruno, I. J., Chisholm, J. A., Edgington, P. R., McCabe, P., Pidcock, E., Rodriguez-Monge, L., Taylor, R., van de Streek, J. & Wood, P. A. (2008). *J. Appl. Cryst.* **41**, 466–470.
- Majumder, M., Buckton, G., Rawlinson-Malone, C. F., Williams, A. C., Spillman, M. J., Pidcock, E. & Shankland, K. (2013). *CrystEngComm*, **15**, 4041–4044.
- Marder, T. C. (2004). Unpublished single-crystal data.
- Minitab (2010). *Minitab 17 Statistical Software*, <http://www.minitab.com>.
- Nichols, G. & Frampton, C. S. (1998). *J. Pharm. Sci.* **87**, 684–693.
- Nishibori, E., Ogura, T., Aoyagi, S. & Sakata, M. (2008). *J. Appl. Cryst.* **41**, 292–301.
- Noguchi, S., Fujiki, S., Iwao, Y., Miura, K. & Itai, S. (2012). *Acta Cryst. E* **68**, o667–o668.
- Noguchi, S., Miura, K., Fujiki, S., Iwao, Y. & Itai, S. (2012). *Acta Cryst. C* **68**, o41–o44.
- Nowell, H., Attfield, J. P., Cole, J. C., Cox, P. J., Shankland, K., Maginn, S. J. & Motherwell, W. D. S. (2002). *New J. Chem.* **26**, 469–472.
- Pérez Cáceres, L., López-Ibáñez, M. & Stützel, T. (2014). *Evolutionary Computation in Combinatorial Optimisation*, Lecture Notes in Computer Science, Vol. 8600, edited by C. Blum & G. Ochoa, pp. 37–48. Berlin, Heidelberg: Springer.
- Post, M. L. & Horn, A. S. (1977). *Acta Cryst. B* **33**, 2590–2595.
- R Core Team (2011). *R: a Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, <http://www.R-project.org/>.
- Rohlíček, J., Maixner, J., Pažout, R., Husák, M., Cibulková, J. & Kratochvíl, B. (2010). *Acta Cryst. E* **66**, o821.
- Rukiah, M. & Al-Ktaifani, M. (2011). *Acta Cryst. C* **67**, o166–o170.
- Rukiah, M. & Assaad, T. (2010). *Acta Cryst. C* **66**, o475–o478.
- Rukiah, M., Lefebvre, J., Hernandez, O., van Beek, W. & Serpelloni, M. (2004). *J. Appl. Cryst.* **37**, 766–772.
- Schmidt, M. U., Ermrich, M. & Dinnebier, R. E. (2005). *Acta Cryst. B* **61**, 37–45.
- Sergeev, G. B., Sergeev, B. M., Morosov, Y. N. & Chernyshev, V. V. (2010). *Acta Cryst. E* **66**, o2623.
- Shankland, K. (2005). *IUCr Commission on Crystallographic Computing Newsletter*, No. 5, pp. 92–102.
- Shankland, N., David, W. I. F., Shankland, K., Kennedy, A. R., Frampton, C. S. & Florence, A. J. (2001). *Chem. Commun.* pp. 2204–2205.
- Shankland, N., Love, S. W., Watson, D. G., Knight, K. S., Shankland, K. & David, W. I. F. (1996). *J. Chem. Soc. Faraday Trans.* **92**, 4555–4559.
- Shankland, K., McBride, L., David, W. I. F., Shankland, N. & Steele, G. (2002). *J. Appl. Cryst.* **35**, 443–454.
- Shankland, K., Spillman, M. J., Kabova, E. A., Edgeley, D. S. & Shankland, N. (2013). *Acta Cryst. C* **69**, 1251–1259.
- Shanmuga Sundara Raj, S., Fun, H.-K., Zhang, J., Xiong, R.-G. & You, X.-Z. (2000). *Acta Cryst. C* **56**, e274–e275.
- Shin, H. S., Song, H., Kim, E. & Chung, K. B. (1995). *Bull. Korean Chem. Soc.* **16**, 912–915.
- Smrčok, Ľ., Jorík, V., Scholtzová, E. & Milata, V. (2007). *Acta Cryst. B* **63**, 477–484.
- Sorrenti, M., Catenacci, L., Cruickshank, D. L. & Caira, M. R. (2013). *J. Pharm. Sci.* **102**, 3596–3603.
- Spek, A. L. (2003). *J. Appl. Cryst.* **36**, 7–13.
- Spillman, M. J., Shankland, K., Williams, A. C. & Cole, J. C. (2015). *J. Appl. Cryst.* **48**, 2033–2039.
- Steiner, T. (2000). *Acta Cryst. C* **56**, 876–877.
- Streek, J., van de, Brüning, J., Ivashevskaya, S. N., Ermrich, M., Paulus, E. F., Bolte, M. & Schmidt, M. U. (2009). *Acta Cryst. B* **65**, 200–211.
- Vallcorba, O., Latorre, S., Alcobé, X., Miravittles, C. & Rius, J. (2011). *Acta Cryst. C* **67**, o425–o427.
- Yang, B., Guang, L., Säntti, T. & Plosila, J. (2012). *Learning and Intelligent Optimization*, Lecture Notes in Computer Science, Vol. 7219, edited by Y. Hamadi & M. Schoenauer, pp. 307–322. Berlin, Heidelberg: Springer.
- Yatsenko, A. V., Chernyshev, V. V., Paseshnichenko, K. A. & Schenk, H. (2001). *Acta Cryst. C* **57**, 295–297.