



profiling

*for success*



abstract

Reasoning Tests

# User's Guide



numerical

Angus S McDonald



verbal



**Please note:**

This publication may not be resold, rented, lent, leased, exchanged, given or otherwise disposed of to third parties. Neither the purchaser nor any individual user employed by or otherwise contracted to the purchaser may act as agent, distribution channel or stockist for this publication. No portion of this publication may be reproduced stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise without the prior written permission of Team Focus Limited. No part of this Booklet is reproducible under any photocopying licence scheme. This Booklet is excluded from the Copyright Licensing Agency Rapid Clearance Service (CLARCS). Team Focus Limited is the worldwide Publisher and Distributor of the Profiling for Success range of instruments\*.

\*Other questionnaires in the Profiling for Success series include:

**Personality and Style:** TDI, LSI, 15FQ+, Emotional Intelligence (EIQ<sup>3D</sup>)

**Teams and Relationships:** MTRi, ITPQ, Relational Health Audit (RHA), 360 models

**Motivation and Interests:** Career Interest Inventory (CII), Value-based Indicator of Motivation (VbIM)

**Capability:** Verbal, Numerical, Abstract Reasoning, Decision Analysis, Memory and Attention

This booklet and all Profiling for Success materials are published by Team Focus Limited of Heritage House, 13 Bridge Street, Maidenhead, Berkshire, SL6 8DS [www.teamfocus.co.uk](http://www.teamfocus.co.uk) email: [info@teamfocus.co.uk](mailto:info@teamfocus.co.uk)

# Profiling for Success: Reasoning Tests

## User's Guide v1.2

Angus S. McDonald

### Contents

Introduction	1
Section One: Using Reasoning Tests for selection and development	3
Why use Reasoning Tests?	3
Reasons for developing the PfS-Reasoning Tests	6
Section Two: Selecting and administering the Reasoning Tests	9
Introduction	9
Selecting appropriate tests	9
Using the PfS- Reasoning Tests as locator tests	11
Administering paper-based tests	13
Overview of administration	13
Planning the test session	13
Materials	14
The test session	15
Administering computer-based tests	17
Supervised assessment	17
Unsupervised assessment	18
Technical requirements for computer-based tests	20
Section Three: Scoring and review of test results	23
Overview of scoring and test scores	23
Qualitative analysis of results	24
Scoring paper-based tests	25
Scoring computer-based tests	26
Using the online report generator with paper-based tests	27
Review of test results	27
Communicating test results	27
Conducting a review session	29

Section Four: Development of the Reasoning Tests	33
Test formats	33
Verbal reasoning format	34
Numerical reasoning format	34
Abstract reasoning format	35
Item writing	36
Pre-trialling item reviews	36
Trialling	36
Item analysis	37
Section Five: Technical information	39
Introduction	39
Reliability	39
The concept of reliability	39
Reliability statistics	40
Standard error of difference	44
Bias	46
A commentary on interpreting bias data	52
Validity	53
Face validity	53
Content validity	54
Construct validity	54
Criterion validity	59
References	63
Appendix One: Explanations of practice questions	65
Appendix Two: Sample test reports	71
Appendix Three: Norm tables	79
Introduction to the norm tables	79
General norms for closed tests	80
Descriptions of norms for open tests	93
Additional norms for closed tests	97
Descriptions of additional norms for open tests	108

## List of tables

- Table 1: Correspondence between the PfS-Reasoning Tests and level of ability as indicated by the level of educational attainment
- Table 2: Appropriate Verbal, Numerical and Abstract test levels for locator test percentile scores
- Table 3: Summary descriptions for combinations of speed of working and accuracy
- Table 4: Score bands used in summary reports and their relationship to T-scores and percentiles
- Table 5: Timings and number of items in each of the PfS- Reasoning Tests
- Table 6: Mean, SD, sample size, number of items, internal consistency and SEM for the PfS-Reasoning Tests
- Table 7: Mean and SD for first time and retest candidates, and test-retest reliabilities for bespoke versions of the PfS-Reasoning Tests
- Table 8: Difficulty levels for the closed PfS-Reasoning Tests and parallel form reliability
- Table 9: Mean raw scores and standard deviations for males and females on the PfS-Reasoning Tests
- Table 10: Mean raw scores and standard deviations for 'whites' and 'non-whites' on the PfS-Reasoning Tests
- Table 11: Mean test scores and effect sizes for different ethnic groups based on the open Level 2 PfS-Reasoning Tests Reasoning Tests
- Table 12: Associations between raw PfS-Reasoning Tests and respondents age
- Table 13: Intercorrelations of the PfS-Reasoning Tests
- Table 14: Associations between PfS-Reasoning Tests and the GMAT
- Table 15: Associations between PfS Abstract Tests and GMA Abstract form A
- Table 16: Inter-correlations between the Verbal, Numerical and Abstract Reasoning Tests and existing reasoning tests
- Table 17: Associations between GCSE English, maths and science grades and PfS-Reasoning Tests

Table 18: The association between UCAS points, degree class and PfS-Reasoning Tests

**List of figures**

Figure 1: The predictive validity and popularity of different assessment methods

Figure 3: PfS-Reasoning Test levels and summary information

Figure 3: The normal distribution curve, Z-score, T-score and percentile scales

## Introduction

The Profiling for Success-Reasoning Tests (PfS-Reasoning Tests) offer a flexible approach to the assessment of reasoning abilities for selection and development purposes. The tests cover three areas of reasoning abilities:

- Verbal – The ability to understand written information and determine what follows logically from the information.
- Numerical – The ability to use numerical information to solve problems.
- Abstract – The ability to identify patterns in abstract shapes and generate and test hypotheses.

As the benefits of psychometric assessments are increasingly recognised and test usage grows, new ways of assessing abilities are needed. The PfS-Reasoning Tests meet these needs by offering both paper- and computer-based assessments that can be used with a wide range of ability groups.

The key features and benefits of the PfS-Reasoning Tests are:

- Flexible delivery options – the paper- and computer-based (online) tests allow for traditional individual or group administration, or remote assessment. Through remote assessment it is possible to include test data earlier on in the assessment process. For example, including test information in the first sift alongside application forms or CVs gives more information on which to base decisions, so potentially enhancing the accuracy of decisions and increasing the efficiency of the selection process.
- ‘Open’ and ‘closed’ test versions – closed versions of each of the tests are available for use under supervised conditions where the identity of tests takers can be closely monitored. Open access versions are also available for use in situations where remote, unsupervised administration is appropriate. These different versions have been developed to meet the increasing need to test candidates remotely, largely as a result in the growth of internet assessment, and the demand for the use of tests for guidance and other development purposes, as well as the more established approach of supervised assessment.
- Common formats across a wide ability range – the Verbal, Numerical and Abstract Reasoning Tests span a wide range of ability levels, from school leavers to experienced managers, using common test formats. If necessary, the appropriate test level can be identified by administering a test as a ‘locator’ among a group of current employees. This process is readily achieved through the use of online tests and guidance is given on how to do this in the User’s Guide.
- Detailed reports and analysis – separate computer-generated reports are available for test users and test takers. For test takers these reports give raw and

standardised test scores and an analysis of speed and accuracy, linked to a narrative suggesting areas for consideration and development. Test users' reports present full test data and an analysis of speed and accuracy linked to interview prompts. Summary versions of reports for test takers and test users are also available.

This User's Guide provides test users with the information they need to understand, use and interpret the Verbal, Numerical and Abstract Reasoning Tests which make up the PfS-Reasoning Tests. **Section One** summarises research on the importance of reasoning abilities for successful job performance and training, and describes the rationale behind the PfS-Reasoning Tests. Administration of paper- and computer-based versions of the tests is covered in **Section Two**. **Section Three** deals with scoring and feedback. The development of the PfS-Reasoning Tests is described in **Section Four**, and **Section Five** provides technical information on the tests and their functioning. It is recommended that all users should read at least **Sections Two** and **Three** before using any of the tests.

In addition to the information contained in this User's Guide, the test publishers offer consultancy, training and general support in using and interpreting the results from these Reasoning Tests and other assessments. For enquiries and support, please contact Team Focus Ltd on + 44 (0)1628 637338, e-mail [teamfocus@teamfocus.co.uk](mailto:teamfocus@teamfocus.co.uk).

## Section One: Using Reasoning Tests for selection and development

### Why use Reasoning Tests?

The use of reasoning tests for selection and development is well-established in many organisations. Surveys show that usage continues to increase (e.g. CIPD, 2006; Jenkins, 2001), with new organisations discovering the benefits that properly applied psychometrics can bring and established users expanding their use of psychometrics. The use of online tests as part of the selection process has also grown rapidly in recent years, with figures showing a rise from 6% in 2002 to 25% in 2006 (CIPD, 2004; 2006).

When used sensitively, with due regard for both their strengths and limitations, there are many good reasons for using psychometric tests. The most compelling reason for using psychometrics is that they provide accurate information on a person's potential or development needs. All benefits of psychometric assessments ultimately feed into this accuracy, so helping the decision-making or development process. Well-informed decisions, in turn, help organisations to grow and develop. It is now well-established that tests of general mental ability, of which reasoning is a core feature, are an important factor in the decision-making process, as they are the best single predictor of job performance and success on work-related training courses (Schmidt and Hunter, 1998).

To contribute to the decision-making process, psychometric tests have to discriminate among the people who take them. Here, discrimination is the ability to identify real differences between test takers' potential, not the pejorative sense of discrimination where one group is favoured over another for reasons unrelated to true potential.

Changes in the education system, particularly the increasing number of students in further and higher education, have made psychometric tests valuable decision-making tools for employers for three reasons:

- The growth in the number of courses and qualifications makes it difficult to evaluate applicants with very different qualifications.
- The increasing number of students obtaining top grades means that academic qualifications have lost much of their ability to discriminate between people.
- Standards of education vary considerably between institutions and courses. Psychometric tests overcome these variations by providing a 'level playing-field' for people to demonstrate their current ability and potential.

This last point touches on the increasingly important issue of fairness in selection. A very significant reason for using psychometrics is that they can provide a fair assessment of all applicants. To be fair, the abilities assessed by the test must be related to job performance (see page 9) and administration standardised for all test takers (see Section Two). Helping test takers to prepare for the testing session, for

example by sending out the Test Taker's Guide (see page 14) or giving access to other approved practice materials, also helps to give everyone an equal chance to demonstrate their abilities.

Psychometric tests further contribute to effective selection and development decisions by explicitly recognising the potential for error in test scores. All assessments (e.g. educational qualifications, ratings from assessment centres or interviews) are subject to error, but this error is rarely acknowledged (see pages 36 and 37 for further discussion of test error). Recognising that test scores contain a degree of error and making this explicit, allows the band of error to be taken into account when making decisions based on test scores.

The relationship between test scores and subsequent job performance or success on training courses has been touched on above. To be defensible as a selection method, links between test scores and subsequent job or training performance have to be established. When this link is established, a test or other selection method is said to have 'validity' or to be 'fit for the purpose'. Showing a test has validity is also important as it is the basis for showing a selection process to be defensible from a legal perspective.

Early research on the links between test scores and subsequent job performance produced mixed results, often due to the limitations of the research itself. More rigorous research methods have since identified a considerable relationship between performance on the family of tests represented by the PfS-Reasoning Tests and job performance (e.g. Bertua, Anderson and Salgado, 2005; Schmidt and Hunter, 1998).

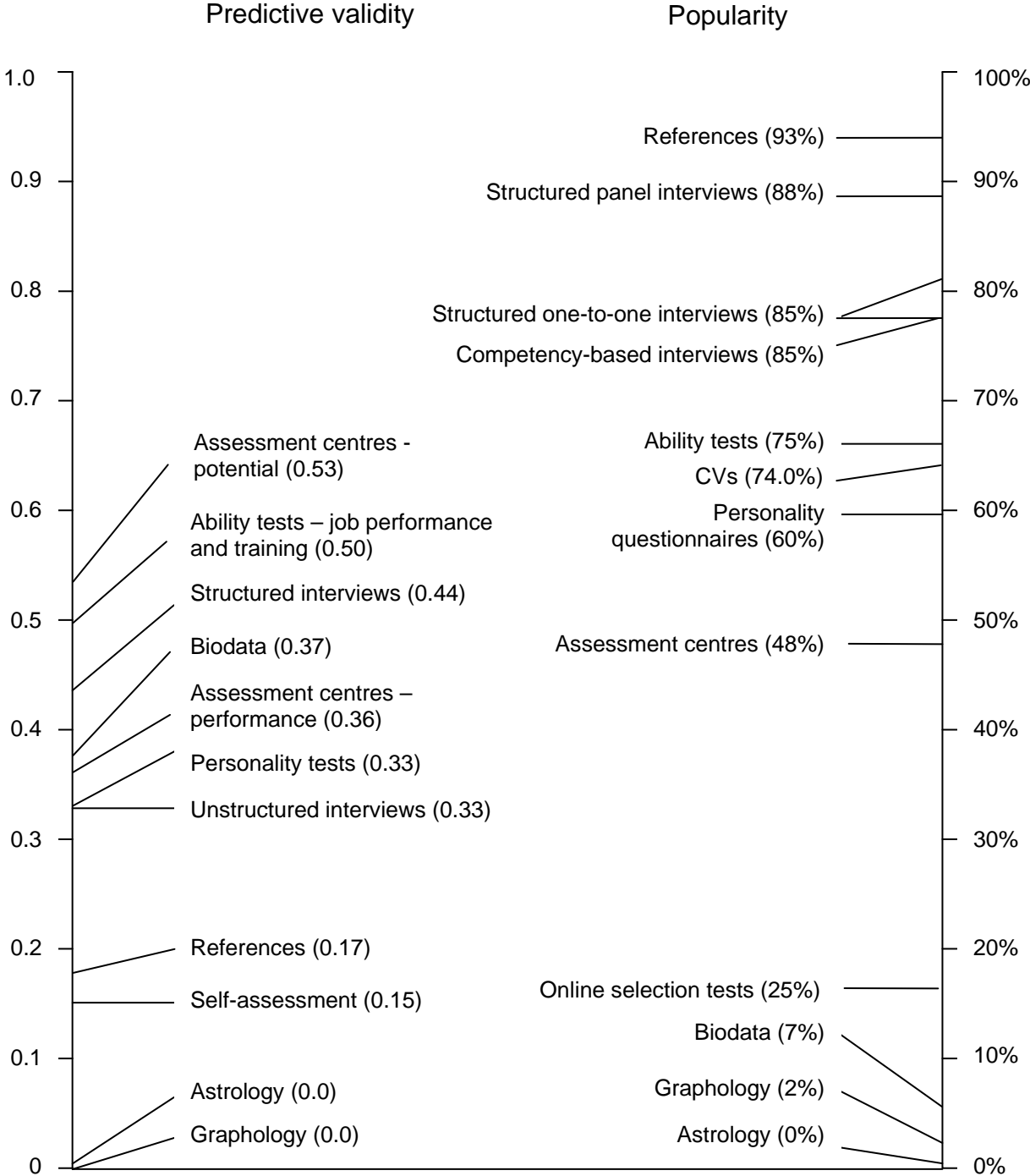
**Figure 1** summarises the findings from a number of sources on the predictive validity and popularity of a variety of assessment methods. From this it can be seen that ability tests are very good predictors of job performance and job related training success, and are one of the most frequently used assessment methods after interviews and references.

In a meta-analysis of validity data, Schmidt and Hunter (1988) showed tests of general mental ability to have a predictive validity of 0.51. Recent work using validity studies from the UK produced figures of 0.48 for the relationship between general mental ability and job performance, and 0.50 for the relationship with job related training success (Bertua *et al*, 2005). Although assessment centres have a slightly higher predictive validity, reasoning tests are often administered as part of assessment centres. The incremental validity of assessment centres once ability tests have been allowed for is quite modest, with recent estimates suggesting that, at best, they add no more than 0.1 to the correlation with job performance (Schmidt and Hunter, 1998).

A further finding of note from Bertua *et al* (2005) was the relationship between predictive validity and different occupational groups. Tests of general mental ability showed higher validities for the 'managerial' and 'professional' categories, indicating their importance for the prediction of more complex, cognitively demanding roles. As

the authors note, this finding contradicts the assumption held by some that ability tests have less validity for more senior appointments.

Figure 1: The predictive validity and popularity of different assessment methods



Notes: figures for predictive validity taken from Beruta, Anderson and Salgado (2005), Gaugler, Rosenthal, Thornton and Bentson (1987), Hunter and Hunter (1982), McDaniel, Whetzel, Schmidt and Maurer (1994), Reilly and Chao (1982), Robertson and Kinder (1993), Schmidt and Hunter (1998). Figures for popularity based on British organisations and taken from CIPD (2000; 2006) and Shackleton and Newell (1991) and indicate use by organisations for at least some appointments.

## Reasons for developing the PfS-Reasoning Tests

The PfS-Reasoning Tests were primarily developed to:

- meet the demand for new materials, due to the increase in psychometric testing;
- offer test users the advantages of computer-based assessments;
- give users the options of 'open' and 'closed' versions of tests assessing the same constructs to address the need for both supervised and unsupervised assessment; and
- deal with common issues with psychometric testing – issues that have been identified from extensive experience of using and training people to use psychometric assessments. These include having sufficient choice to select tests of appropriate level of difficulty according to the ability range of the individuals being tested and to have a consistent format across these levels in order to reduce the learning curve (sometimes needed since tests can have very different formats) for both administrator and testee.

The advantages of computerised testing have been recognised for some time, particularly in the areas of administration, scoring and generating reports (e.g. Kline, 2000). Within the fields of occupational, clinical and educational assessment, computer-based testing is now widely accepted. Test developers are starting to go beyond computerised versions of pencil-and-paper assessments and explore how technology can be used to create innovative and engaging new forms of assessment.

When developing the PfS-Reasoning Tests, one of the goals was to make the potential benefits of testing more accessible. Using the internet as a method of delivery means that the psychometric assessments can be used more flexibly and with a wider range of people. With computer-based testing it is also easier to make tests such as the PfS-Reasoning Tests visually appealing. This is important when assessing people who may have lower levels of motivation for completing the tests, as it makes the testing experience different from traditional educational assessments.

The PfS-Reasoning Tests also meet the need for administrators to be able to set up and monitor the assessment process, to have control of the data and how it is used, and to generate informative reports from the test results. Through the PfS online system, administrators can control which tests are made available to test takers, which norm groups the results are compared to, what types of report are generated and who receives the reports. Security is guaranteed by the use of passwords. Computerised testing makes scoring fast and accurate.

The output from the PfS-Reasoning Tests provides some added value information which goes beyond the usual raw and standardised test scores. They provide an analysis of speed and accuracy (see **Section Three**) to enable the interpreter to consider potential reasons for low or high scores which may have to do with strategy rather than just ability. These statistics are combined into reports that give

development suggestions and interview prompts, through which reviewers can explore with test takers the meaning and implications of their results. These analysis and reporting facilities mean that all test takers can receive valuable, personalised feedback, regardless of the outcome of the test results. This makes the PfS-Reasoning Tests truly developmental. By using the data entry and scoring facilities of the PfS online assessment system, users of the paper-based tests can also benefit from the features of the automated tests reports (see pages 24 and 25).

Two related challenges often faced by test users are:

- selecting the appropriate ability level of psychometric tests; and
- familiarising themselves with the formats of different tests.

Both of these issues are addressed by the PfS-Reasoning Tests.

It is not usually possible to adequately cover a wide ability range with a single test. Tests provide maximum information when the mean score is close to the middle of the possible score range. A single test capable of assessing accurately across the full ability range would take too long to administer to be practical in organisational settings and would be frustrating for many test takers: able test takers would become bored with many simple questions and less able test takers frustrated at the number of questions they found too difficult. To deal with the issue of ability, there are four levels of the Verbal, Numerical and Abstract Reasoning Tests, spanning school-leavers to people with postgraduate qualifications and considerable professional experience.

Each of the three tests uses a common format. This means that once users have familiarised themselves with one level of a test, they will be equally familiar with all levels of both the closed and open versions. Users of the PfS-Reasoning Tests therefore no longer have to become familiar with different test formats for different populations – simplifying the administration process. The same formats are also used for scoring and reporting, reducing the possibility of errors and making the interpretation and review of test results easier.

Test users often see identifying the appropriate test level as a major challenge, particularly when the test is to be used by a diverse group of people. The PfS-Reasoning Tests address this issue through suggesting how the tests can be used as ‘locator’ tests. By administering one of the tests to an existing group of employees, the results can be used to determine which of the four test levels is appropriate for the position in question. The use and interpretation of locator tests is simplified if the computer-based versions are used. Guidance on how to use and interpret locator tests is given on pages 10 to 12.

In areas such as graduate and managerial selection and development, the use of psychometrics is well-established. As more organisations use psychometrics there is a risk that the tests become over-exposed, with applicants in some cases taking the same test more than once, so giving them an unfair advantage over others.

All new tests offer a short-term solution to the problem of over-exposure, though this has become an increasingly important issue with the advent of unsupervised testing over the internet. The PfS-Reasoning Tests have also been developed with the goal of addressing this in the long-term. The open and closed versions of the Verbal, Numerical and Abstract Reasoning Tests have been developed to give users confidence in the security of the closed tests whilst retaining the option of unsupervised internet assessment using the open versions. Further parallel versions of the Verbal, Numerical and Abstract Reasoning Tests are already under development and there is also the option for bespoke assessments consisting of unique series of items to be developed for clients on request. The PfS-Reasoning Tests therefore offer organisations the opportunity to avoid the problems associated with the over-exposure of tests.

The range of PfS-Reasoning Tests is illustrated below, showing which are available as open and closed tests, in paper and online formats, and the number of items and timing for each.

*(diagram to illustrate open and closed, with N no. of items and timing)*

To summarise the PfS-Reasoning Tests available:

- Levels 1 to 4 closed tests cover the areas of Verbal, Numerical and Abstract Reasoning and are intended to be used for secure testing situations which are either supervised or where they are administered online to known test takers.
- Each level contains a unique set of items and levels are broadly tied to educational stages: Level 1 for test takers in the last years of compulsory education (years 10 and 11), Level 2 for those in further education, Level 3 for undergraduates and Level 4 for postgraduates and experienced professionals.
- Levels 1 and 2 of the open tests are intended for use under less secure conditions (e.g. during an initial sift where results are collected remotely).
- As with the closed tests, each level contains a unique set of items. Level 1 of the open tests covers the same ability range as Levels 1 and 2 of the closed tests and Level 2 of the open tests the same range as Levels 3 and 4 of the closed tests.
- The Combined Reasoning Test consists of items from the Level 1 open tests. As such it is intended for use under less secure conditions, particularly initial sifts and career development or guidance in younger test takers.

## **Section Two: Selecting and administering the Reasoning Tests**

### **Introduction**

For any test to play a valuable role in the decision-making process, it has to be matched to the abilities and competencies required by the job role. The first part of this section provides an overview of how to identify appropriate tests and introduces the facilities in the PfS-Reasoning Tests series that allow the most suitable level of each test to be selected.

Good administration, whether the tests are being taken in pencil-and-paper format or using a computer, is the key to achieving reliable and valid test results. When administering the test in person, a well-defined procedure is to be followed. However, computer administration offers test takers the opportunity to complete tests in their own time, at a location of their choosing, without an administrator being present. Under these conditions the administration procedure may not be as closely controlled, but it is still possible for clear guidelines to be established. The second part of this section outlines the procedure for supervised test administration and goes on to offer guidelines for organisations on how to develop procedures for unsupervised testing.

### **Selecting appropriate tests**

The information provided by the PfS-Reasoning Tests should be valuable in the decision-making process. To make sure this is the case, the abilities being assessed by the tests must relate to core job competencies. The starting point for any selection or development process must be a detailed job analysis focussing on the competencies and personal characteristics that employees need in order to perform successfully. As job roles and organisational structures become ever more fluid, identifying and assessing the competencies needed for those who work in these changing environments can also help organisations plan for future development and growth.

It is important to remember that reasoning tests can provide valuable information, but are rarely sufficient on their own. Tests should be seen as only one part of an overall assessment package. As with any form of assessment, both their strengths and weaknesses need to be acknowledged. Through drawing on the strengths of a variety of assessment methods and carefully integrating the information from them, it is possible to reach far more valid and defensible decisions that are also more likely to be viewed as fair within the framework of employment law.

In order to provide maximum information about individuals, it is important that the correct level of each Reasoning Test is selected. If tests are not at the correct level for the group in question, their ability to differentiate between people is lowered and they may have a de-motivating effect on those who take them. It is important to recognise that selecting more difficult tests will not result in a raising of standards within an organisation. Tests give most information when scores are spread around the mid-point of the distribution; if they are too easy or too hard, scores will be more

bunched together so making it difficult to reliably differentiate between the test takers. The availability of appropriate norm groups is another factor in determining test selection and also indicates for which ability levels or groups tests are suitable.

Currently, there are four levels of each of closed Reasoning Tests and two levels of the open Reasoning Tests (referred to in the online PfS assessment system as 'Reasoning Skills Tests' to differentiate them from the closed tests). In addition, there is also the Combined Reasoning Test which includes verbal, numerical and abstract items in a single test. For this Combined Reasoning Test there is just one level. Each level has been developed to correspond to a broad ability band, as shown in **Table 1**. These bands should be considered as a starting point for test selection.

Reasoning Test level	Reasoning Skills Test level	Approximate educational level of the norm group
Level 1	Level 1 and Combined Reasoning Skills Test	This covers the top 95% of the population and is broadly representative of the general population.
Level 2		This covers the top 60% of the population and is broadly representative of people who have studied for A/AS Levels, GNVQ Advanced, NVQ Level 3 and professional qualifications below degree level
Level 3	Level 2	This covers the top 40% of the population and is broadly representative of the population who study for a degree at a British University or for the BTEC Higher National Diploma/Certificate, NVQ Level 4 and other professional qualifications at degree level
Level 4		This covers the top 10% of the population and is broadly representative of the population who have a postgraduate qualification, NVQ Level 5 and other professional qualifications above degree level

*Table 1: Correspondence between the PfS-Reasoning Tests and level of ability as indicated by the level of educational attainment*

When using the tests with well-defined groups, such as A-level students or those just about to graduate, **Table 1** should be adequate for appropriate test selection. Deciding on the appropriate level is more difficult when the tests are being used for a less homogenous group, particularly when some people may have considerable work experience but limited academic qualifications. When the most suitable test level is not immediately apparent, users may consider identifying the appropriate level by using a 'locator' test. A description of how to use a locator test is given below.

## Using the PfS- Reasoning Tests as locator tests

To identify the appropriate level of the PfS-Reasoning Tests, it is possible to use one of the tests as a locator test. Either the paper- or computer-based tests can be used in this way, but it is more efficient and flexible to use the computer-based tests. The locator test approach is possible because of the common format of the PfS-Reasoning Tests, and simplified by the time-efficient nature of the tests and computer-based scoring and reporting.

By administering locator tests to current employees, it is possible to establish a mean level of reasoning abilities within specific groups. This information can then be used to select the most appropriate level test for the position in question. It is suggested that Level 2 of the closed PfS-Reasoning Tests are used as the locator tests, as these should not be found too difficult by employees and will give an appropriate indication of which of the four levels is most suitable. If only one of the three test types – Verbal, Numerical or Abstract – will eventually be used, then this one should be the locator test. If two of three tests are being used, it is suggested that the test with the highest level of face validity for the role in question is used so as to get most buy-in and motivation from the volunteers being asked to complete it.

The locator tests provide a method for identifying reasoning abilities in current employees. It is recognised that many employees will have developed their skills in key job areas since being employed through training programmes, job experience or a combination of the two. Although it is not possible to determine the actual extent of skill growth, allowance for this is made through recommended test levels being adjusted slightly downward where borderline scores occur.

As with any test administration, it is important that good practice is followed for the administration of the locator tests if they are to provide valid information. Whilst it is not necessary to conduct a formal group administration, the following stages are recommended:

- Identify the group of current employees to take the locator test. As far as possible, this group should be at the same level as those being selected and working in the same job role.
- Ideally between 10 and 20 people should take the locator test. Asking for volunteers is likely to result in a sample that is more comfortable, and probably capable, with reasoning tests. It is not best practice to make it compulsory for employees to take the locator test. Hence, there are two possible ways in which groups can be identified:
  - A random sample from all suitable employees can be taken.
  - If there is a need to raise the skill level in the specific reasoning area, or if tests are to be used for development purposes, a sample can be taken from employees who are known to perform well in their job roles. However, it is important to ensure that the test level identified is not too high, meaning that

tests will not adequately discriminate between potential employees. Excluding the bottom 20 or 25% of performers in the area assessed by the Reasoning Test may be an appropriate cut-off for identifying a high-performing group.

- Selected employees should be contacted, either by letter or email, requesting their participation in the locator testing. The purpose of the testing should be made clear, as should the confidential and anonymous nature of the test results. Guidelines should be given about taking the tests, for example that they should be taken in a quiet environment, free from disturbances. Test takers should also be told how long they should allow for completing the test. Clear information about the locator test is important to ensure that employees ‘buy in’ to the testing process and are motivated to perform at their best.
- The website address and passwords can be sent to employees at the same time as their participation in the testing is requested. Alternatively, employees can be asked whether they are willing to participate and the website address and passwords sent to those who agree. If paper-based tests are being used, details of when and where the test session will take place should be sent out. Administration procedures should follow the guidance given later in this section.
- If employees are not to receive individual results from the locator test, it is important to acknowledge their contribution to the process and to thank them for their participation.

When sufficient data has been collected, the mean of the raw test scores should be calculated. As mean scores can be affected by individual values that are far from the majority of scores (‘outliers’), data should be visually inspected to check whether there are any extreme values. If there are any scores that are 6 or more raw score points below all others, it is recommended that these are removed before the mean is calculated. Lower scores may be particularly a problem due to the motivation of some test takers. Higher scores should not be removed as these will reflect high levels of ability in the sample.

**Table 2** shows the recommended level of the Verbal, Numerical and Abstract Tests, according to mean locator test score. Note that these scores are based on the percentiles from the Level 2 test, using the norms given on pages 83, 87 and 91 for Verbal, Numerical and Abstract respectively

Locator test percentile score	Recommended PfS test level	
	Reasoning Tests (closed)	Reasoning Skills Tests (open)
1-35	Level 1	Level 1 and Combined Test
36-70	Level 2	
71-90	Level 3	Level 2
91-99	Level 4	

*Table 2: Appropriate Verbal, Numerical and Abstract test levels according to locator test percentile scores*

## Administering paper-based tests

### Overview of administration

For a test session to be fair and to fulfil the purpose for which it was designed, it is important that it is run efficiently and smoothly. The way a test session is delivered can potentially affect the anxiety and performance of the test takers, their impression of the organisation and their motivation to perform well. The aim is for the administrator to be personable, efficient and clear when giving the test instructions.

This part of the User's Guide gives full instructions on how to prepare for administering the PfS-Reasoning Tests. In addition, there is a separate card of Administration Instructions for each test, which sets out the exact procedure to follow for the test session. Administrators are advised to prepare using this User's Guide in conjunction with the Administration Instructions and then, in the test session itself, just to use the Administration Instructions and any personal notes they have made.

For each test, administrators need to familiarise themselves with the:

- Question Booklet
- Answers to the example and practice questions
- Answer Sheet
- Administration Instructions card
- Test Log.

Before administering any of the tests, administrators should take the tests themselves – this is the best way to understand what is required. The procedure set out on the Administration Instructions card should be practised and administrators should make sure that they fully understand the solutions to the example and practice questions (full explanations to the practice questions are given in **Appendix One**).

The PfS-Reasoning Tests can be administered in any order although the most usual is:

- Verbal
- Numerical
- Abstract.

### Planning the test session

The test room needs to be suitably heated and ventilated (with blinds if glaring sunlight is likely to be a problem) for the number of people taking the tests and for the length of the test session. The room should be free from noise and interruption, as any disturbances can affect test takers' performance. There should be space between each test taker's desk so that test takers cannot see others' papers and the administrator can walk around.

If the tests are to be taken as part of an assessment day, remember that performance tends to deteriorate towards the end of a long day. If a number of test sessions are being planned, those who take the tests towards the end of the day may be disadvantaged. It is recommended that test takers can take any two of the Reasoning Tests without a break being needed between them. If all three tests are being administered, there needs to be a break, normally between the second and third tests. A break of at least ten minutes is recommended.

If more than 15 candidates are attending the test session, it is advisable for the administrator to have a colleague to assist with the administration. The efficiency of the session will be improved if a second person can check that test takers have grasped the practice questions and format of the Answer Sheets, and to assist with administration generally. Some preparation is also necessary for this role, particularly familiarisation with the Question Booklets, Answer Sheets and explanations to the practice items.

Test takers should be notified of the date, time and location of the test session and told which test(s) they will be taking. The Test Taker's Guide can be sent out at this stage, to help candidates prepare. The Test Taker's Guide is available online, allowing test takers to be contacted by email, if appropriate. This method may be particularly useful if test takers will be completing the computer-based tests. At this point, it is good practice to inform candidates why they have been asked to take the tests, how the results will be used in the selection procedure, how they will receive feedback about their performance and to explain the organisation's policy on the confidentiality of test results. The Test Taker's Guide can be accessed from the following link to the PfS website:

[www.profilingforsuccess.com/about/documents/test\\_takers\\_guide.pdf](http://www.profilingforsuccess.com/about/documents/test_takers_guide.pdf)

When test takers are notified about the session, it is essential that they are also asked to contact the administrator or other appropriate person, if they have any disabilities that will affect their ability to complete the tests and to specify what accommodation needs to be made for them to complete the tests. Under the Disability Discrimination Act (1995; 2005), test users are obliged to make changes to assessment procedures so that people with disabilities are not disadvantaged at any stage of the selection process. By obtaining information about any special needs well in advance of the test session, organisations can make the necessary adaptations to the testing session and have time to seek further advice if necessary. Further information on assessing people with disabilities can be found on the PfS website as:

[www.profilingforsuccess.com/about/documents/Assessing\\_People\\_with\\_Disabilities.pdf](http://www.profilingforsuccess.com/about/documents/Assessing_People_with_Disabilities.pdf)

## **Materials**

Before the testing session, ensure that there are the correct number of Question Booklets and Answer Sheets. Question Booklets should be checked to make sure

that they have not been marked. Marks should be erased if possible, or replacement books obtained. The Test Log has been developed to help administrators prepare for the testing session – it contains a checklist of the materials needed and other arrangements that have to be made. It also allows administrators to record the room layout, any unusual occurrences during the test session and to summarise the test scores of a group of test takers. It is a useful document to keep for later review sessions or if any challenges are made to the test results or decisions that the results feed into.

Each test taker needs:

- a Question Booklet
- an Answer Sheet
- two ball-point pens or pencils (pencils need to be sharp to clearly mark the carbonated answer sheet)
- two sheets of paper for rough working
- a candidate ID number (if applicable).

The administrator needs:

- a copy of the appropriate Question Booklet and Answer Sheet
- the appropriate Administration Instructions card
- a Test Log
- spare pens/pencils
- spare rough paper
- two stopwatches or watches with a second hand
- explanations to the practice questions if not fully familiar with them.

There is space for test takers to record personal information on the Answer Sheets. Not all of this information may be needed, so administrators should make sure they are clear about what information is required and ask test takers to complete only what is necessary.

### **The test session**

A notice to the effect of 'Testing in progress – Do not disturb' should be displayed on the door of the test room. Ensure that chairs and desks are correctly positioned. Place two pens or pencils (these need to be sharp to clearly mark the carbonated Answer Sheet), two sheets of paper for rough working, and ID numbers (if applicable) on each test taker's desk. Do not issue the Question Booklets and Answer Sheets at this stage.

If ID numbers are being used but have not already been allocated to test takers, allocate these outside the test room, then ask test takers to enter the room and find the corresponding desk. Otherwise, invite test takers into the test room and direct them where to sit.

#### *Stage 1: Informal introduction*

When all test takers are seated, the administrator should give the informal introduction to the test session. This needs to be prepared in advance to include the points given below, but should be delivered informally, in the administrator's own words. The aim here is to explain clearly to the test takers what to expect and to give them some background information about the tests and why they are being used. This will help to reduce anxiety levels and create a calm test setting. The administrator should aim for a relaxed, personable, efficient tone, beginning by thanking the test takers for attending.

The important points to include in the informal introduction are:

- Introduce the administrator (and any colleagues, if appropriate) giving their position in the company.
- The programme for the test session including: the timing, which tests will be taken, how long each test will last and the timing of any breaks (use a flipchart to show the programme if it is at all complex).
- Why the organisation is using the tests, who will see the results and how these will be used in the selection process. Explain what will happen to the test results and how they will be recorded and stored, emphasising confidentiality and accessibility in accordance with the Data Protection Act.
- Check comfort levels and whether anyone needs the cloakroom, as test takers will be asked not to leave the room once the tests have begun.
- Explain how test takers will receive feedback about their performance.
- Tell test takers that they will be given full instructions before each test, will be able to see examples, try practice questions and ask questions before the test begins, to make sure they fully understand what they have to do. Confirm that all tests will be timed.
- Ask the test takers if they have any questions so far, and address these.

At the end of the informal introductory talk, test takers should be told that from this point the tests will be administered according to a set procedure and that the instructions will be read from a card, to ensure that all candidates receive exactly the same instructions. The administrator should now turn to point 4 on the appropriate Administration Instructions card and follow the exact procedure and wording given.

### *Stage 2: Formal testing procedure*

Begin the formal testing procedure at point 4 on the relevant Administration Instructions card. It is important to follow the given procedure and wording exactly, to ensure that the instructions are the same and therefore fair and consistently administered to all test takers.

On each Administration Instructions card, the text in the shaded boxes should be read out verbatim to test takers. The text outside of the shaded boxes contains instructions for the administrator.

Use the Test Log to note the number of Question Booklets and Answer Sheets distributed and collected in, to ensure that none go astray. The start and finish time of each test should also be recorded on the Test Log. There is also room on the Test Log to record anything that occurs during the test, relating to individuals (e.g. the need for replacement pens or to leave the test room) or to the group as a whole (e.g. fire alarm or other disturbance). This information can be important later – for example, when comparing the performance of groups from different test sessions, or if an individual queries the basis of his or her selection or other decision based on test performance.

At the end of the test, collect in the Question Booklets and Answer Sheets while the test takers are still seated, ensuring while doing this that each test-taker has entered any required biographical details on the answer sheet and have indicated which level of the test they have taken. If several tests are being administered, replace any pens/pencils and rough paper, if necessary. Start the procedure for the next test from point 4 on the relevant Administration Instructions card. At the end of the session, thank test takers for attending and explain what they should do next.

## Administering computer-based tests

Computer-based testing offers users far greater flexibility than paper-based tests. It also benefits from automated scoring and the ability to produce full reports almost instantly. Procedures for administering computer-based testing, particularly testing over the internet, are not as well-established as for paper-based testing. This part of the User's Guide discusses some of the options for computer-based testing. It does not set out to prescribe a process, but introduces the issues that need to be considered and makes some recommendations, so that users can formulate their own policies in this area. Administering computer-based tests under supervised and unsupervised conditions will now be considered. The technical requirements for the computer-based tests are also described.

### **Supervised assessment**

Computer-based tests can be used as an alternative to paper-based tests. Here, test takers, either as individuals or in groups, complete the tests under supervised conditions as they would paper-based tests. The formal test instructions, example and practice items are given on-screen and so do not need to be read from the Administration Instructions card. An appropriate approach to test administration in this situation would be as follows:

- Check that the room and computers are set up appropriately.
- Invite test takers into the testing room and direct them where to sit.
- Ask test takers not to touch the computers until they are told to do so.
- Give the informal introduction as for paper-based tests (see page 15), but tell the test takers that they will be taking the test on computer.
- At the end of the informal introduction, ask if there are any questions.
- Direct test takers to the PfS website and follow the appropriate link to take a test, then give them the Client code, Access code and Password to enter when

prompted. Alternatively, prior to the beginning of the testing session, ensure that the PfS website has already been accessed on each computer and the entry codes entered in order that the PfS assessment facility is already displayed on screen when candidates take their places at their computers.

- Tell test takers that the computer will prompt them to enter their personal information before giving them the test instructions and practice and example items.
- Test takers should be allowed to work through the instructions at their own pace and begin the test when they are ready.
- Explain that if they have any questions or experience any difficulties during the test, they should raise their hand.

Test takers will finish the tests at slightly different times using this approach, as not everyone will work through the instructions at the same pace. If this approach is taken, administrators should decide whether to ask test takers to remain seated until everyone completes the test or whether they can leave the room when they have finished. This is likely to depend on the number of people being tested and the room set-up (i.e. how easily people can leave the room without disturbing others).

Alternatively, test takers can be asked to work through the instructions, practice and example items, and then wait until everyone is ready to begin. When everyone is ready, the administrator should ask test takers to start. Everyone will finish the testing session at the same time if this approach is used, thus eliminating the possibility of test takers who have been slower to work through the instructions being disturbed by others leaving the room.

If two of the Reasoning Tests are being taken, test takers can be instructed to move on to the second test when they have completed the first. As with the paper-based tests, if all three tests are being used, it is recommended that test takers are allowed a break between the second and third tests.

Finally, it should be noted that the tests which will be displayed on the screen when test-takers enter the PfS assessment area on the PfS web site will depend on the 'Access Code' which has been used to log in to the system. Administrators should therefore ensure that they have set up an Access Code which includes only the appropriate tests and test levels which they wish to be presented. A discussion of access codes is beyond the scope of this manual, though detailed information will be provided by Team Focus to users of the PfS online assessment system.

### **Unsupervised assessment**

The internet offers the potential to exploit the benefits of testing in new ways, but takes users into the less familiar territory of unsupervised assessment. There are many issues with unsupervised assessment: access to technology, fairness and the authenticity of test results being paramount. Despite the need to address these issues, the benefits of internet-based testing are many. Particularly notable are its efficiency and the opportunity to gather additional information to feed into the early stages of the decision-making process.

When planning an unsupervised testing session, administrators need to consider the target group and their likely access to technology. Certain groups (e.g. university students or those already working for an organisation) may have greater access to the necessary technology than others (e.g. people returning to work). Where it is anticipated that a number of potential test takers may not have access to the necessary technology, it may be advisable not to use internet testing unless other appropriate arrangements can be made. For example, it may be possible to direct test takers to places such as libraries, careers centres or an organisation's regional offices where they can take the PfS-Reasoning Tests under appropriate conditions.

Access to the necessary technology is also related to issues of fairness. If completing internet-based assessments is made a compulsory part of an application process, this may bias the process against those who do not have easy access to the necessary technology. In some cases it could also constitute deliberate discrimination and so be unlawful. Although many organisations use online application procedures, alternatives to these should be put in place (e.g. a paper-based test session available on request). Organisations may have to accept that, in some cases, test results will not be available for all applicants.

A major question with any unsupervised testing session concerns the authenticity of results. As the tests are unsupervised, there is no way of telling who has actually completed the tests or whether the intended test taker has received assistance. If the PfS Reasoning Tests are being used for development purposes or careers guidance, authenticity should be less of an issue. It is during selection that issues around authenticity are most critical.

One significant advantage of internet-based testing, as mentioned above, is that psychometric tests can be used early in a selection procedure, possibly at the same time application forms are completed. If used as part of a selection decision, it is essential to be confident that the test results are indeed the work of the applicant.

Ensuring the validity of test results requires that test takers are monitored during the test session. This removes many of the advantages of internet-based testing, so it is important to encourage honesty in test takers. One way in which this can be done is to position the tests as offering potential applicants valid feedback on their abilities and the demands of the job. This would imply on the one hand, suggesting to low scorers that the job may not be well matched to their abilities, and so would be unsatisfying for them and, on the other hand, confirming to higher scorers that they appear to have the necessary basic abilities required by the job. If test scores are used to make decisions at an early stage of an application process, it may be prudent to give them a lower weighting than normal and to set lower standards of performance.

The validity of test scores is more of an issue with high scorers. One approach to dissuade people from obtaining assistance with the tests is to view them as a 'taster' to the next stage of selection where further testing will take place under more controlled conditions. If test takers know that they will have to take a similar test

under supervised conditions if they proceed to the next stage of the selection process, they may be less inclined to seek assistance with the unsupervised tests. In these circumstances it may be appropriate to initially use the open versions of the Reasoning Tests, then follow these up with the closed versions under supervised conditions if it is deemed necessary to verify results.

All the issues discussed above need to be considered when undertaking unsupervised, internet assessment. Despite this, in many ways, the actual test procedure is not that different from supervised administration. The main stages of the test process remain the same, although as it is not possible to give an informal introduction to the test session, the initial contact with test takers is very important. The contact letter, email or telephone conversation should include:

- why they are being asked to take the tests.
- what tests they have to take.
- how the results will be used.
- how they will receive feedback on their test results and who will have access to them.
- the hardware/software requirements of the tests.
- appropriate conditions for taking the tests (how long they should allow, the need for a quiet room, free from disturbances).
- how to access the testing site (website address and passwords).
- when the tests should be completed.
- either a copy of, or web link to, the Test Taker's Guide, recommending that this is used to help prepare for taking the tests.
- what will happen when the tests have been completed.
- the details of who should be contacted in case of queries or difficulties.

Particularly important under unsupervised test conditions will be the information on why the tests are being used. As discussed above, positioning the tests as providing applicants with an insight into their own suitability for the job can help to encourage honesty and acceptance of the remote testing experience when used for selection. If applicants who proceed to the next stage will have to take further tests, this should also be stated, again to encourage honesty.

Once test results have been received, an opportunity should be made for test takers to discuss their results (see **Section Three**).

### **Technical requirements for computer-based tests**

If internet testing is being considered, the issue of access to technology needs to be addressed. Although the majority of people now have access to computers, it should not be assumed that this is the case for everyone. It also needs to be recognised that conditions should be conducive to completing a timed test; some computers that are accessible to the public may be in noisy environments and where test takers are liable to disruption.

To make the PfS-Reasoning Tests widely accessible, the system has been designed to make minimal demands on technology. The system will work on any internet-ready computer. The preferred browser is Internet Explorer version 5 or later with Adobe Flash® version 5 or later installed. The minimum screen resolution needed is 800 x 600 though a resolution of 1024 by 768 is recommended. Virtually all modern desktop computers and most modern laptop computers will meet the specifications needed to run the tests. Tests are accessed over the internet. As the whole test is downloaded before the test begins, timing for the test is unaffected by the speed of the internet connection.

It is not necessary for the internet connection to be maintained once a test has been downloaded. However, the internet connection does have to be active when the test results are submitted. Information about the need for test takers to be actively connected to the internet for their test results to be recorded is displayed at the end of the test.



## Section Three: Scoring and review of test results

### Overview of scoring and test scores

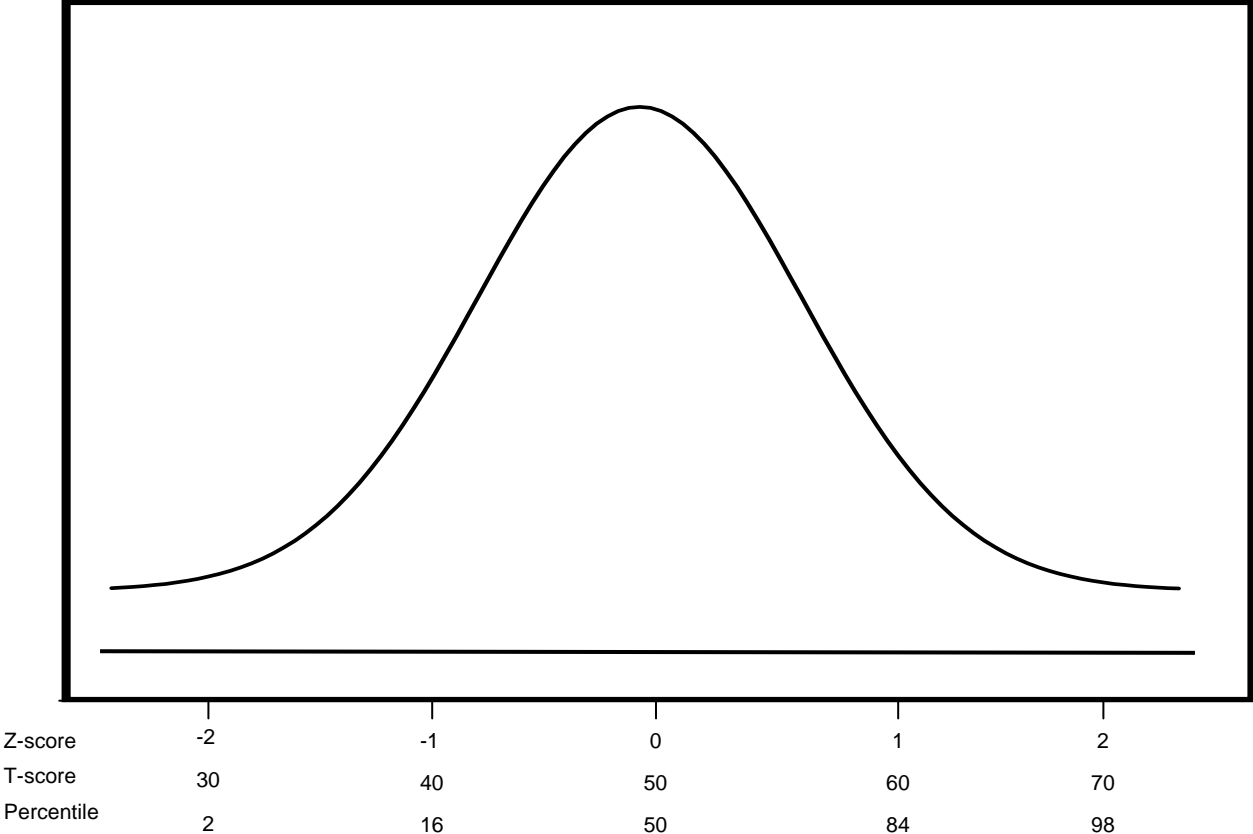
The primary purpose of testing is to obtain a test score or mark which says something about the test taker's ability. To understand any test score however, it needs to be put into context. For example, if a person has answered 26 out of a possible 30 questions correctly, this appears initially to be a good score. But if another group of other test takers all score between 27 and 30 on the same test, is 26 still a good score?

The purpose of this example is to highlight that simple test scores cannot be considered 'good' or 'poor', without knowing more about how people generally perform on the test. Test scores are put into context by comparing them with the scores of a large group of people who have previously taken the test. This group is known as the 'norm group' and the tables that allow individual scores to be compared to those from the norm group are called 'norm tables'. The norm tables for the PfS-Reasoning Tests are in **Appendix Three**. The types of scores given in the PfS-Reasoning Tests norm tables are described below.

- Raw score – The raw score is the number of marks a test taker achieves on a test. For the Verbal, Numerical and Abstract Reasoning Tests, one mark is given for each question that is answered correctly. Therefore, the raw score is the number of questions answered correctly.
- Percentile – Percentiles describe the proportion of the norm group a test taker has scored the same as or better than. For example, a percentile of 65 means that the person has scored as well as or better than 65 percent of the norm group. Because percentiles are quite easy to understand, they can be particularly useful when communicating information to test takers or people who are unfamiliar with psychometric testing.
- T-score – T-scores are a transformation of the raw scores onto a scale which is approximately normally distributed (that is, a bell-shaped distribution with no long tails). This transformation is necessary as raw score distributions are often skewed, with more scores towards the higher or lower end of the distribution. T-scores have a mean of 50 and a standard deviation (SD; an indication of the spread of scores) of 10. The main advantage of using a scaled score such as T-scores, is that they allow performance on different tests to be directly compared.
- T-score confidence band – All test scores contain a degree of error, as no test can give a perfect indication of a person's ability. This error can be quantified and described as a range within which a person's 'true score' is likely to fall. The norm tables give 68% and 80% confidence bands for each T-score. These confidence bands indicate the range in T-scores between which it is 68% or 80% certain that a person's true score will lie. For a more detailed discussion of test error and confidence bands see the section on Reliability (pages 41 - 47).

The relationship between percentiles, T-scores and the normal distribution curve are shown in **Figure 3**.

*Figure 3: The normal distribution curve, Z-score, T-score and percentile scales*



**Qualitative analysis of results**

In addition to providing a single test score, the PFS-Reasoning Tests combine two further test statistics to produce a more qualitative assessment of each test taker’s performance; the number of questions the test taker has attempted and the proportion of questions attempted that have been answered correctly. Both of these values can be compared to norm tables in the same way as the raw score and are classified as being ‘above average’, ‘average’ or ‘below average’. This process results in a three by three matrix, describing the test taker’s speed of working and accuracy as shown in **Table 3** overleaf.

		Accuracy (proportion of questions attempted answered correctly)		
		Below average	Average	Above average
Speed (number of questions attempted)	Below average	Slow and inaccurate	Slow and moderately accurate	Slow and accurate
	Average	Average speed and inaccurate	Average speed and accuracy	Average speed and accurate
	Above average	Fast and inaccurate	Fast and moderately accurate	Fast and accurate

*Table 3: Summary descriptions for combinations of speed of working and accuracy*

The analysis of speed and accuracy has been developed to help both test takers and users gain a fuller understanding of the Reasoning Test results. The nine summaries of test performance given in **Table 3** have been expanded into more detailed descriptions of performance, interview prompts and development suggestions. These descriptions are included in the full reports generated by computer-based tests or from the test scoring facility on the Profiling for Success website. Summary reports provide an overview the test taker's speed and accuracy but do not include full interview prompts or development suggestions.

When using the extended analyses in the full versions of the reports, it needs to be recognised that these reports offer a range of possible reasons for a person's performance. These are offered as ideas and prompts that can be used during a review session to explore test performance in more detail. Every effort has been made to make these reports comprehensive, although they should not be seen as exhaustive. Further, the reports attempt to reflect the test taker's ability, approach to the test and mental processes, but may be limited in some cases as the descriptions can be extrapolated only from the test taker's responses.

## Scoring paper-based tests

Answer Sheets for the paper-based tests are made up of two sheets of paper; the top sheet where test takers mark their answers and the bottom sheet which contains the scoring key. As the top sheet is carbonated, the test taker's biographical information and answers are transferred on to the bottom sheet with the scoring key. The steps that should be followed when scoring Answer Sheets are set out below.

1. Before beginning to score an Answer Sheet, check that the test taker has completed all the necessary biographical details and that they have indicated which level of the test they have taken.

2. On the right hand side of the Answer Sheet there is a perforated strip. Tear off this strip and then use a pencil or ruler to separate the top and bottom pages of the answer sheet.
3. Check that the test taker's personal details and answers to the questions have transferred clearly on to the bottom part of the Answer Sheet.
4. Count up and tick the number of times the responses given by the test taker correspond to the correct answers indicated on the bottom part of the Answer Sheet. As each correct answer is worth one mark, the total number of ticks is their 'raw score'. Enter the raw score in the box marked 'Raw score'.
5. Count up the number of questions to which the test taker has given an incorrect or ambiguous response, and add this to their raw score. This gives the number of questions that the test taker has attempted. Enter the number of questions they have given a response to in the box marked 'Number of questions attempted'.
6. Use the appropriate norm table to look up the percentile, T-score and confidence bands that correspond to the T-score. These should be entered in the appropriate boxes on the answer sheet.
7. On the reverse of the bottom part of the Answer Sheet test takers may have recorded comments about the test and the test session. This information should be available to the person conducting the review session, as it can provide useful information to discuss during the review.

Sometimes test takers make ambiguous marks on answer sheets. The following guidelines for resolving ambiguities were applied during the development of norms for the paper-based Reasoning Tests. These should be followed to ensure the validity of normative data.

- If more than one answer has been indicated to a question, and all but one answer is clearly crossed out, count this as the intended answer and score against the scoring key.
- If more than one answer has been indicated to a question, score as incorrect.
- If all answers have been crossed out, score as incorrect.
- Test takers may miss out a question but forget to leave a blank space for the question on their answer sheet. This is most apparent when a series of answers are incorrect according to the scoring key, but each indicates the correct answer for the following question. If a series of four or more answers indicate the correct answer to the following questions, it is possible that an answer has been missed out. In such cases, appropriate adjustments should be made and the questions treated as correct.

## Scoring computer-based tests

Computer-based tests are automatically scored when the answers are submitted at the end of the test. From the scored data a report including the raw score, percentile,

T-score and confidence bands are automatically created for each test taker. An extended analysis, as described on pages 22 and 23, is also included if the full version of the report is requested. This report is sent to the email address entered by the administrator during the set-up stage of the testing process.

When setting up an assessment in the PfS 'Client Area', there is also an option for users to request reports for the test taker. Test taker's reports are versions of the main reports in a format that can be given directly to the test taker. As with the administrator's reports, full or abbreviated versions of these reports are available. If test taker's reports have been requested, these will also be sent to the email address entered by the test-taker when logging in to the PfS system.

Samples of full reports and a summary reports can be seen in **Appendix Two**.

## Using the online report generator with paper-based tests

Users of the paper-based tests can also make use of the online report generator that is a standard part of the computer-based tests. The report generator requires users to be set up as clients of the Profiling for Success system, which is accessed via the internet at the following address:

[www.profilingforsuccess.com](http://www.profilingforsuccess.com)

The test system will ask users to enter their Client Code and Password. Test data can then be entered through the 'direct data entry' screens. Reports will be generated on the submission of the data. For more information on this system or to set up an online PfS account please contact Team Focus.

## Review of test results

Whenever psychometric tests are administered it is good practice to review the results with the test taker. The exact format of this review will depend on the purpose of assessment and how the results are to be used. Practical considerations, such as the number of people being tested and the opportunities to meet with test takers if assessment is being conducted over the internet, will also affect how the review of test results is delivered. However, regardless of the specific situation, test takers should always be given the option to discuss their results and to raise any questions they have. The following sections provide guidance on how test results can be communicated to test takers and how to conduct a review session.

### Communicating test results

There are three main ways in which results can be communicated to test takers. These are outlined below, along with some of the considerations around each method.

- *Face-to-face review session.* The preferred method of communicating test results is to do so in person (guidance on conducting a personal review session is given

below). Face-to-face reviews have the advantages of encouraging openness and honesty, allowing reviewers greater insight into the test taker's reactions to the results and so opportunities to respond to these, and generally encourage greater interaction. The results from a Verbal, Numerical and Abstract Reasoning Test can be covered in a single review session lasting between 10 and 15 minutes, so these review sessions do not have to be time-consuming. These can be scheduled to follow testing sessions or interviews, to avoid difficulties in arranging subsequent meetings for the reviews. The option of reviewing results almost immediately after tests have been completed is possible due to the rapid scoring and report generating facilities of the PfS-Reasoning Tests, particularly when the computer-based tests are used.

- *Telephone review.* When there is no personal contact with the test taker (for example, when initial screening has been conducted over the internet and some candidates have not progressed to the next stage of assessment), telephone review sessions can be conducted. A mutually convenient time needs to be arranged between the reviewer and test taker, to ensure that both have sufficient time, free from interruptions, for the review to be conducted fully. A particular limitation of this approach is that reviewers do not have access to non-verbal cues, which can be valuable in gauging a test taker's reactions during face-to-face reviews. Under these conditions, reviewers need to be particularly aware of emotional reactions in what test takers say and may need to prompt more around how the test taker is feeling about their results than when conducting face-to-face reviews.
- *Written review.* Giving test takers purely written information on their test performance is the least preferable way of communicating results. This method essentially gives 'feedback' (test results being delivered to the test taker) as there are very limited opportunities for exploring the meaning and implications of the test results. Whenever this method is used, it is important that test takers are given the name and telephone number of a person they can contact to discuss their results.

The PfS-Reasoning Tests can produce reports for both the reviewers and test takers. The test takers' reports (see **Appendix Two**), either in their full or summary versions, are suitable for giving directly to test takers as care has been taken in developing these reports to ensure that the language used is accessible and generally positive. Test takers' reports give raw scores and percentiles in written and graphical form, and include personal development suggestions. The score bands used in the administrator's reports and their relationship to T-scores and percentiles are shown in **Table 4**.

Score band used in summary report	T-score band	Percentile band
Low	36 and below	<1-10
Below average	41-37	11-30
Average	42-58	31-69
Above average	59-63	70-89
High	64 and above	90-99

*Table 4: Score bands used in administrator's reports and their relationship to T-scores and percentiles*

Reviewers need to consider if they want to use the test takers' reports and, if so, when they will be given to test takers. It can be useful to give the reports to test takers before the review session, allowing time for the reports to be read and test takers to think about issues they may want to raise in the review session. In principle, reports can be given to test takers during the review session or at the end of it. However, if reports are given during the review session, test takers may be distracted by the reports unless time is allowed for them to be read.

Alternatively, reviewers can use information gained through the review session to edit and tailor reports before giving them to test takers. This approach may be particularly useful in developmental contexts, when personal development suggestions and action plans can be included in the final report given to test takers.

If a report is to be edited after a review, it is suggested that the following process is used:

1. The report and any associated graphics files are saved from the email to a folder on the computer.
2. The report file can be opened by a word processing package such as Microsoft Word<sup>®</sup>. To do this it may be necessary to select the option to allow html file types to be viewed and opened.
3. The report can then be edited as a normal Word<sup>®</sup> document and saved in its original html format.

### **Conducting a review session**

The purpose of a review session, whether conducted face-to-face or via the telephone, is to ensure that the test taker clearly understands the meaning of their results, is satisfied with the assessment experience and to explore possible implications of the results. To reach these goals it is important that the review session is seen as a chance for information to be given and received by both the test taker and the reviewer, not simply for the reviewer to provide the test scores. For this process to be successful, it is vital that all reviewers have received appropriate training.

General guidelines for conducting review sessions are given below. These guidelines should be seen as identifying the main points that need to be covered and

giving suggestions about the structure of the review session and appropriate questioning strategies. They do not set out to provide a set formula that must be followed. Although the guidelines below are written with face-to-face reviews in mind, they are also applicable to telephone reviews.

- As with test administration, good preparation is essential for review sessions. A suitable room, free from disturbances, should be identified. Reviewers should familiarise themselves with the individual's test results, what the test measures and how this relates to the job role, and any other relevant biographical information. Technical language should not be used during the review session, so it is useful for reviewers to prepare a simple description of what each test measures. For example, a Numerical Reasoning Test may be better described as 'an opportunity to show how comfortable you are with using numbers and numerical information to solve problems'. Reports should be sent to test takers in good time if these are being given out before the review session.
- The review session should begin with the reviewer introducing themselves and providing a brief overview of the purpose of the review session. Useful information to provide includes the approximate length of the session, issues around confidentiality and what will happen to the test results.
- Both parties need to agree on what they want to get out of the session, such as information, consequences of test performance or a way forward.
- To encourage a balanced discussion from the outset, the test taker should be brought into the review session as early as possible. This can be done through asking the test taker about their experiences of the tests immediately after the brief introduction (e.g. "How did you find the reasoning tests?" or "Tell me about your experience of taking the tests?"). Throughout the review session open questions should be used wherever possible, as this will encourage the test taker to provide more information and make the review more balanced. In a balanced review session there should be equal contributions from both the reviewer and the test taker.
- If the tests were completed some time before, a reminder of these and how they fit into the selection or development process may need to be given at this stage.
- At this point it is also appropriate to explain how test results are interpreted with reference to a norm group. It is generally best to avoid the term 'norm group' as this may not be understood by all test takers and for some may imply 'normal' performance. A preferable phrase is 'comparison group', which conveys the process of comparing individual test scores to those from a wider group of people, and is more readily understood.
- The next stage involves discussion of the actual test results. It is preferable to let the test taker lead the order in which the tests are reviewed, rather than going through the tests in order. The review process can be started through questions such as "Which test did you prefer and why?" or "Which test did you find most

challenging?”. Once a test has been identified, the reviewer can give the test taker their score or can ask them to estimate their own performance on the test, for example “In relation to the comparison group (describe comparison group) how do you feel you performed on the (appropriate test)?”.

- It is preferable to describe test scores in terms of percentiles, though it needs to be clearly communicated that percentiles refer to the proportion of the comparison group who the test taker scored as good as or better than, and not the percentage of questions they answered correctly. It may also be informative to explore the number of questions that test taker attempted and number answered correctly as this, in conjunction with the text around speed and accuracy, can be used to explore the way in which the test taker approached the test.
- Once the test taker’s performance on each test has been established, their reactions to the result and its implications need to be explored. For example, questions such as “How do you feel about your result on this test?” can be used to assess emotional reaction and “What implications do you think the test results may have on your application?” or “How might your result on this test influence your choice of career?” can be used to explore the implications of test results. Although reviewers often perceive low scores as more challenging to discuss, it is important that low test scores are not ‘glossed over’ or dismissed. Questions such as “How far do you think the result is a fair reflection of your ability in this area?” can be very valuable. Often test takers have a reasonable insight into their abilities and low scores in some areas may not necessarily be a great source of concern; test takers often find it quite reassuring to know that they have performed at an ‘average’ level.
- If the computer-generated test user’s reports are being used, these contain interview prompts to explore test performance and the implications of this. As these reports combine information on speed and accuracy, they offer prompts that are specifically tailored to the individual’s performance and how they could improve their performance. Because of this they can be particularly valuable when using the tests for development or when the reviewer has limited experience of working with this type of assessment.
- The final stage of the review process is to ask the test taker to summarise what has been discussed, to ensure clear understanding. Summaries can take the form of a brief review of the test results that highlight any strengths and weaknesses that have been identified. The implications of the test results and any development plans should also be summarised, if these have been discussed. To check that the test taker has understood what has been discussed, it can be valuable to get them to summarise what they see as the main points to have emerged from the review session, rather than this being provided by the reviewer. The reviewer should explain the next stage in the selection or development process and what will happen to the results, and inform the test taker about confidentiality. Finally, the test taker should be offered the opportunity to ask any outstanding questions and then thanked for attending the review session.

It is good practice for individual organisations to develop policies around the review of test results, as with other aspects of testing. Such policies should cover the organisation's general policy on test reviews, how reviews are conducted, confidentiality and storage of information. It is important for organisations to develop their own policies, as these will help ensure consistency of approach and application over time, and will also guard against issues of fairness and discrimination. Whilst policies may draw on the guidelines given above, ultimately reviewers should develop their own style, with which they feel comfortable, within these frameworks.

## Section Four: Development of the Reasoning Tests

### Test formats

The development of the Verbal, Numerical and Abstract Reasoning Tests involved a number of stages. The purpose of the first stage was to define as clearly as possible, the final format of the tests. By understanding how the final tests would look and function, the test development team identified the main determinants of the item formats. The key aspects affecting item formats were identified as:

- Ability range – The tests should be suitable for a wide range of abilities; from people with average GCSE passes or equivalent, up to postgraduates and those with considerable professional experience. It was therefore necessary to identify item formats that could support questions across this ability range.
- Computer technology – From the outset it was decided that the tests should be primarily computer-based, but that many users would still want pencil-and-paper versions to be available. Item formats that could be used in both mediums were therefore needed. The test development team wanted to exploit the advantages of computer-based technology, but also recognised that this technology needed to be widely available. While the internet was seen as offering the greatest flexibility for testing, it was necessary to consider issues such as use of graphics, download times and the possible unreliability of internet connections in the design of the tests.
- Test length – Test users may often want to incorporate a broad range of assessments into their selection and development procedures, but also need to consider the time available for testing. With this in mind, the target time for each test was set at 15 to 20 minutes, depending on test level. Item formats therefore needed to be efficient (e.g. by minimising the time spent on unnecessary reading or calculations) but also needed to remain sufficiently contextualised to reflect real-life problems. The tests were therefore designed to be problem-solving tasks, presenting test takers with information and questions based on the information. To make the items 'efficient' a number of questions were related to each piece of information, or 'stem', and the length of the stems was carefully controlled.
- Answer format – To allow for the scoring of the tests to be achieved quickly and reliably, a multiple-choice format was used for each of the tests. Whilst open-response items can offer a rich source of information, this is offset by the difficulties in scoring open-response tests reliably, particularly when scoring tests by computer. The time needed for scoring and resolving ambiguities in open-response tests was also seen as unacceptable to those who use high volumes of tests.

## **Verbal reasoning format**

The Verbal Reasoning Tests consist of passages of information, with each passage being followed by a number of statements. Test takers have to judge whether each of the statements is true or false on the basis of the information in the passage, or whether there is insufficient information in the passage to determine whether the statement is true or false. In the latter case, the correct answer option is 'can't tell'.

As test takers come to the testing situation with different experiences and knowledge, the instructions state that responses to the statements should be based only on the information contained in the passages, not on any existing information that test takers have. These instructions also reflect the situation faced by many employees who have to make decisions on the basis of information presented to them. In these circumstances decision-makers are often not experts in the particular area and have to assume the information is correct, even if they do not know this for certain.

The passages of information in the Verbal Reasoning Tests cover a broad range of subjects. As far as possible, these have been selected so that they do not reflect particular occupational areas. Passages were also written to cover both emotionally neutral areas and areas in which people may hold quite strong opinions or have emotional involvement. Again, this was seen to make the Verbal Reasoning Test a valid analogy of decision-making processes, where individuals have to reason logically with both emotionally neutral and personally involving material.

Each statement has three possible answer options – true, false and can't tell – giving test takers a one-in-three or 33% chance of guessing the answer correctly. Guessing is most likely to become a factor when tests are highly speeded. The quite generous time limits and the 'not reached' figures, suggest guessing is unlikely to be a major factor for the Verbal Reasoning Tests. The proportion of true, false and can't tell answers was balanced in both the trial and final versions of the Verbal Reasoning Tests. The same answer option is never the correct answer for more than three consecutive statements.

## **Numerical reasoning format**

The Numerical Reasoning Tests present test takers with numerical information and ask them to solve problems using that information. Some of the harder questions introduce additional information which also has to be used to solve the problem. Test takers have to select the correct answer from the list of options given with each question.

Numerical items require only basic mathematical knowledge to solve them. All mathematical operations used are covered in the GCSE (Key Stage 4) mathematics syllabus, with problems reflecting how numerical information may be used in work-based contexts. Areas covered include: basic mathematical operations (+, -, x, ÷), fractions, decimals, ratios, time, powers, area, volume, weight, angles, money, approximations and basic algebra. The tests also include information presented in a variety of formats, again to reflect the skills need to extract appropriate information

from a range of sources. Formats for presentation include: text, tables, bar graphs, pie charts and plans.

Each question in the numerical test is followed by five possible answer options, giving test takers just a one-in-five or 20% chance of obtaining a correct answer through guessing. The distractors (incorrect answer options) were developed to reflect the kinds of errors typically made when performing the calculations needed for each problem. The answer option 'can't tell' is included as the last option for some problems. This is included to assess test takers' ability to recognise when they have insufficient information to solve a problem. As with the Verbal Reasoning Tests, the same answer option is never the correct answer for more than three consecutive statements.

### **Abstract reasoning format**

The Abstract Reasoning Tests are based around a categorisation task. Test takers are shown two sets of shapes, labelled 'Set A' and 'Set B'. All the shapes in Set A share a common feature or features, as do the shapes in Set B. Test takers have to identify the theme linking the shapes in each set and then decide whether further shapes belong to Set A, Set B or neither set.

The abstract classification task is based on Bongard problems (Bongard, 1972). Bongard problems were originally developed to test the ability of computer-based pattern recognition programs. In their original form these problems consisted of two sets, each containing six shapes. Computer programs had to identify the common feature(s) of the shapes in each set, but they were not required to classify further shapes.

A development of this task was chosen for the Abstract Reasoning Test as it requires a more holistic, inductive approach to problem-solving and hypothesis-generation than abstract problems involving sequences of shapes or progressions. People operating at high levels are often required to focus on different levels of detail, and to switch between these rapidly (e.g. understanding budget details and how these relate to longer-term organisational vision). These skills are assessed through the Abstract Reasoning Test, as it requires test takers to see patterns at varying levels of detail and abstraction. The lower level of the abstract test can be a particularly valuable tool for spotting potential in young people or those with less formal education, as it has minimal reliance on educational attainment and language.

Test takers are required to identify whether each shape belongs to Set A, Set B or neither set. This gives three possible answer options, meaning test takers have a one-in-three chance of guessing answers correctly. As with the other tests, the proportion of items to which each option is the correct answer has been balanced. The same answer option is never the correct answer for more than four consecutive shapes.

## Item writing

The test items were written by a team of people, who all had extensive experience of occupational psychology or using assessments in selection and development contexts. Detailed briefing notes were assembled for item writers, outlining the nature of the tests, the specific details of the items for each test type and giving example items. Prior to writing test items, all item writers attended a workshop which introduced them to the process of item writing and covered areas of good practice, particularly in relation to bias. This was followed by a practical session involving item writing and group review of the items produced.

After attending the workshop, item writers initially submitted a limited number of items to the test development team for review and feedback. Only after these initial items were considered satisfactory were item writers given the go-ahead to develop the test items. Throughout the item writing process, feedback was continually given to item writers to ensure standards were maintained and that adequate coverage of each area was achieved.

## Pre-trialling item reviews

Before trialling of the items took place, they went through a number of review stages. As they were submitted, items were reviewed by the test development team to ensure they met the test specifications, were accurate and unambiguous, and free from bias. All items were further reviewed by an occupational psychologist who was not involved with the test development.

Following the internal item reviews, items were sent to external specialists. Verbal items were reviewed to check for clarity of language and reasoning, Numerical items were reviewed for language and mathematical accuracy (including plausibility of distractors), and Abstract items were checked for accuracy and ambiguity in solutions. Verbal items, Numerical items and instructions for all three Reasoning Tests were also reviewed by an educational psychologist who specialises in language and cultural issues to ensure they were accessible and free from gender and ethnic bias.

## Trialling

The majority of the trialling was computer-based, as this allowed the tests to be accessible to a wide range of organisations and individuals and reduced the possibility of errors in data collection and transfer. Computer-based trialling also allowed timing data to be collected for each test item. An analysis of timing data was included in the initial item analyses and contributed to the final selection of items, so helping to make the tests efficient in terms of time.

Each trial test lasted between 30 and 35 minutes. The trial tests were designed to be less speeded than the final tests, so that sufficient data could be collected on items later in the tests. However, a timing element was included to create realistic test conditions for the trial tests, as this was needed for accurate data for analysis and

item selection. The trialling design involved placing common items in adjacent levels of each test to allow linking and the substitution of items between test levels as necessary.

The trialling software collected biographical information on each test taker. Information on age, gender, educational qualifications and ethnicity were collected for all people who took part in the trialling. Further information was also collected from specific samples as appropriate (e.g. current course of study and predicted degree grades for graduate samples). In total, almost 2000 people participated in the trialling of the items for the closed tests between January 2002 and July 2002. Trialling for the open tests was conducted between October 2002 and February 2003. Approximately 3000 people took part in this exercise.

## Item analysis

For trialling, the timed part of the tests lasted between 30 and 35 minutes and each had approximately double the number of items of the final tests. Once sufficient trialling data had been collected, each test went through a series of analyses to identify items that were not functioning satisfactorily.

Item analyses were conducted to identify the facility and discrimination for each item. The facility indicates the proportion of people who attempted the item who answered it correctly, effectively indicating the difficulty of the item. The discrimination is the point biserial correlation between the score on the item (1 for correct, 0 for incorrect) and total test score excluding that item. This statistic indicates the degree to which each item is able to distinguish between people who obtained high overall test scores and those who obtained lower scores.

As each of the tests uses a multiple-choice format, 'distractor' analyses were conducted on the incorrect items. Essentially, these are the same as discrimination analyses, but examine the link between each of the incorrect answer options and total test score. If items are functioning well, people who choose incorrect answers should get lower overall test scores. If this is not the case, items may be ambiguous, leading strong test takers to choose incorrect answer options.

The time taken by test takers to answer each question was also recorded, and the mean times for each of the questions calculated. This analysis identified items that were answered particularly quickly or slowly on average. Timing information was used primarily to ensure that the items selected for the final versions of each test obtained maximum information within the times allowed for each test. This analysis also complemented the item analysis, suggesting where items may be too easy (very rapid response times) or ambiguous (slow response times).

Each of the trial tests was subjected to a bias analysis to ensure that selected items were not found disproportionately easy or hard by different groups of test takers. Comparisons were made between males and females, and 'whites' and 'non-whites'. Items displaying significant levels of bias were excluded, or were included but balanced with an item showing an equal degree of bias for the opposite group.

Following the item analysis, tests were assembled for standardisation. The time allowed and number of items in each of the standardisation tests are given in **Table 5**.

Closed Reasoning Tests				
		Verbal	Numerical	Abstract
Level 1	Time allowed	12	12	10
	Number of items	32	28	50
Level 2	Time allowed	12	12	10
	Number of items	32	28	50
Level 3	Time allowed	15	15	12
	Number of items	40	36	60
Level 4	Time allowed	15	15	12
	Number of items	40	36	60
Open Reasoning Tests				
		Verbal	Numerical	Abstract
Level 1	Time allowed	15	15	12
	Number of items	44	40	70
Combined test	Time allowed	10	10	7
	Number of items	24	20	35
Level 2	Time allowed	20	20	15
	Number of items	60	48	75

*Table 5: Timings and number of items in each of the PfS-Reasoning Tests*

## Section Five: Technical information

### Introduction

This section of the User's Guide provides a detailed account of the technical functioning of the PfS-Reasoning Tests, covering the areas of reliability, bias and validity. The important area of reliability of measurement and the precision of test scores is explored in detail here, although the key reliability statistics – internal consistency, standard error of measurement and standard error of difference – are also summarised with each of the norm tables in **Appendix Three**.

### Reliability

#### The concept of reliability

No test, including those in the PfS-Reasoning Tests series, gives a perfect indication of reasoning ability. Despite rigorous test development and appropriate use and administration of the tests, there will always be some degree of error in any test result. The concept of reliability is concerned with quantifying the amount of error in a test score. If the accuracy of a test score is known, then scores can be used sensitively with due regard for this error. Reliability is also important as it sets the upper limit on validity: a test cannot be more valid than it is reliable.

The need to take error into consideration is important in many situations, but it is vital when tests are being used to make important decisions that affect people's lives (e.g. recruitment and development decisions). Good psychometric tests have the advantage that their error is made explicit. In many other forms of assessment, no recognition of error is made and test scores or results are treated as absolute truths. A good example of this is exam grades or degree classes, which often contain more error than psychometric tests despite there being no acknowledgement of this error.

According to classical test theory, any test score is made up of two components: true score and error score. A person's true score is their hypothetical score on the trait being measured. For the PfS-Reasoning Tests, the true scores refer to a person's Verbal, Numerical or Abstract reasoning ability. However, scores obtained from tests also contain an error component. Error in test scores can come from three sources: the test itself, the person taking the test and the situation in which the test is being taken.

- Test error – Classical test theory assumes tests are made up from a sample of items taken from the universe of all possible items. As with any sample, this will contain a degree of error. As all people taking a test answer the same set of items, test error is systematic error, being the same for each test taker. Providing that adequate content validity has been ensured, test error is less of a concern to test users than individual or situational error.
- Individual error – The individuals who take the tests are a source of random error. Factors such as how the person is feeling, their motivations and attitudes towards

the testing session, and their familiarity with tests and the test format will all affect how they perform, but are not necessarily related to their actual ability. Sending out the Test Taker's Guide is one way to help limit the effect of individual error, as it ensures all test takers have a chance to become familiar with the format of the tests and know how to prepare for the test session.

- Situational error – The actual test session itself is a further source of random error. The guidelines on administration and the standardised instructions aim to make each testing session as similar as possible for all test takers. However, it is not possible to standardise the testing situation completely. The rooms used for testing, environmental conditions, time of day and interaction between the administrator and test takers will all vary between sessions. Each of these factors can influence test performance but are not related to the test taker's true ability.

### **Reliability statistics**

In practice, the reliability of a test is typically assessed in three ways. The first of these is to look at how the test items hang together to form a coherent assessment of the construct under consideration. This 'internal consistency' is found by taking the mean of the correlation between each test item and total test score, excluding that item. Internal consistency is calculated through a formula known as Cronbach's Coefficient Alpha (or Kuder-Richardson 20 (KR20) when test items are dichotomous) and expressed as a statistic that can range from 0 to 1. The closer to 1, the more reliable the test is said to be.

The second way in which reliability is assessed is through looking at how consistent results are over time. This is done through administering the test at one point in time and then again sometime later. The scores from the two administrations are then correlated with each other to give an indication of 'test-retest' reliability. As with internal consistency, the closer the test-retest correlation coefficient is to 1, the more reliable the test is seen to be.

A further way in which reliability can be assessed is through parallel, or alternate, forms of the test. Alternate versions of the same tests can be particularly useful in applied settings, where it may be desirable to administer the same test more than once or to use a less exposed version. Typically, parallel forms are administered back-to-back and the results from the two are correlated, as when assessing test-retest reliability.

Each of the statistics described above provides an index of reliability, but does not directly indicate the degree of error in a given test score. The standard error of measurement (SEM) provides a way of quantifying the error in a test score, indicating the range within which a person's true score is likely to fall. The SEM is derived from the following formula:

$$SEM = SD\sqrt{1-r}$$

where the *SD* is the standard deviation of the test in raw score units and *r* is the reliability (in this case internal consistency) of the test.

The SEM is used to create 'confidence bands' around test scores. It is known that a person's true score is likely to fall within one SEM either side of their observed score, 68% of the time. This range of scores around an observed score is known as a 'confidence band'. By multiplying the SEM by 1.28, 1.65 or 2, the confidence band can be increased to 80%, 90% or 95%. Using these values it is possible to be 80%, 90% or 95% certain that a person's true score will fall within the confidence band. In the norm tables given in **Appendix Three**, 68% and 80% confidence bands around the T-scores are given. The following sections present evidence on the internal consistency, test-retest reliability and parallel form reliability.

### ***Internal consistency***

**Table 6** shows the descriptive statistics, internal consistency and SEM for each of the Reasoning Tests. A number of factors can affect the reliability and the SEM statistics. A brief discussion of the two main factors follows to allow tests users to understand the statistics given in **Table 6** more fully.

The length of a test is an important determinant of its reliability. Classical test theory assumes that any test is made up of a sample of items from the domain being assessed. As with any sample, the results from it should be more accurate as the sample becomes larger. Hence, there is a trade-off between reliability and practicality: high reliability is desirable, but if a test takes a long time to complete, very few people will choose to use it.

It is possible to construct highly reliable tests of manageable length, by developing them carefully. The rigorous development process for the PfS-Reasoning Tests is described in **Section Four**. As development was done using computer-based tests this also allowed timing data on each test item to be gathered during the trialling stage, meaning that time-efficient items were selected for the final tests. This has resulted in the timed part of the PfS-Reasoning Tests – between 10 and 15 minutes for the closed tests and 15 and 20 minutes for the open tests – being less than many similar tests of equivalent or even lower reliability.

Another factor that affects reliability is the time limit allowed for the test. When tests are highly speeded, reliability estimates tend to become inflated. The item analyses indicate that the time limits allowed for each of the tests to be fairly generous, with the 'not reached' figures being similar to comparable tests. Reliability estimates are therefore unlikely to be unduly affected by the timing of the tests.

Test and level		Mean	SD	Sample size	Number of items	Internal consistency	SEM
<b>Closed Reasoning Tests</b>							
Verbal	1	16.62	5.73	210	32	0.90	1.81
	2	16.32	5.18	303	32	0.80	2.32
	3	24.10	6.07	1322	40	0.86	2.27
	4	25.45	6.27	1131	40	0.87	2.26
<b>Open Reasoning Tests</b>							
Numerical	1	19.30	4.64	250	28	0.93	1.23
	2	14.95	4.74	337	28	0.84	1.90
	3	18.04	5.69	1609	36	0.87	2.05
	4	16.24	6.50	1510	36	0.89	2.16
Abstract	1	28.51	7.82	156	50	0.93	2.07
	2	20.80	8.24	242	50	0.87	2.97
	3	31.20	11.18	860	60	0.92	3.16
	4	30.35	10.41	881	60	0.91	3.12
Verbal	1	14.90	12.37	1010	44	0.92	3.50
	2	29.61	10.32	24072	60	0.91	3.10
	C*	13.75	4.70	763	24	0.84	1.88
Numerical	1	14.45	10.76	1356	40	0.92	3.04
	2	18.31	6.48	37241	48	0.85	2.51
	C	12.35	4.15	763	20	0.86	1.55
Abstract	1	39.04	13.47	515	70	0.95	3.01
	2	33.69	11.67	13.61	75	0.92	3.30
	C	16.83	5.99	763	35	0.85	2.32

\* Combined Reasoning Test

*Table 6: Mean, SD, sample size, number of items, internal consistency and SEM for the PfS-Reasoning Tests*

### **Test-retest reliability**

Evidence of the test-retest reliability for the PfS-Reasoning Tests has been obtained from a client who requested bespoke versions of the Verbal, Numerical and Abstract

tests for their selection process. These tests consist of items from the Levels 2, 3 and 4 closed tests plus other items taken from the Reasoning Test item bank. The tests taken by candidates at this organisation are computer-based and taken under supervised conditions. The organisation's policy allows candidates to re-apply after a period of time if they are initially unsuccessful, so giving a subset of applicants who have two sets of Reasoning Test data.

The sample for the test-retest analysis consisted of 169 candidates who first completed the tests April 2003 and May 2005, and completed them for the second time (retest) between July 2003 and November 2005. One hundred and thirty seven (81.1%) were male and 32 (18.9%) were female. Mean age at time of first testing was 21.4 years (SD=3.8). The mean length of time between first taking the tests and retesting was 38.7 weeks (SD=25.3 weeks), with a range from 2 days to 121 weeks. For the majority of candidates retesting occurred between 10 and 40 weeks after first taking the tests.

	First time (n=169)		Retest (n=169)		Difference	Test-retest correlation	Number of test items
	Mean	SD	Mean	SD			
Verbal	26.7	5.2	29.0	5.0	2.3	0.73	40
Numerical	18.1	3.8	19.6	3.6	1.5	0.71	36
Abstract	42.3	9.5	49.6	9.9	7.3	0.67	70

*Table 7: Mean and SD for first time and retest candidates, and test-retest reliabilities for bespoke versions of the PfS-Reasoning Tests*

Test test-retest correlation coefficients showed that each of the tests had adequate reliability, particularly considering the extended time between testing for many in the sample – almost half a year on average. It should also be noted that the test-retest correlations in **Table 7** are likely to underestimate the true correlations, as these have not been corrected for measurement error.

The data in **Table 7** also gives an indication in likely change in scores on retest. Mean scores on all three tests increased on retesting, although the standard deviation remained relatively constant. An indication of the magnitude of score change can be obtained by looking at the absolute change in mean test score as a proportion of the test's standard deviation (taken from all 5294 candidates from this organisation). From these calculations values of 0.43, 0.38 and 0.71 for the Verbal, Numerical and Abstract tests respectively were obtained, showing modest increases in mean Verbal and Numerical scores and a slightly larger increase in Abstract scores of just under three quarters of a standard deviation. One possibility for the modest changes in score is the extended time period before retesting for some candidates (over 2 years in a few cases). Correlations between the difference between first time and retest scores, however, gave no indication that these scores were associated with time between the two test sessions.

## Standard error of difference

Given that all test scores contain a degree of error, one important question which test users often ask is “Are the scores of two people really different?”. If test scores were free from error, any difference in observed scores would reflect a real difference in the ability being assessed. However, because of error, if two scores are close to each other, there is a chance that they could be reversed if the test takers took the tests again. In other words, the person who obtained the higher score may not continue to obtain the higher score the second time around.

The likelihood of the difference between two test scores reflecting a real difference in the construct being assessed, can be determined with a statistic known as the standard error of difference (SED). The SED indicates how far two test scores need to be apart before the difference can be seen as meaningful. The formula for the SED is:

$$SED = \sqrt{SEM_1^2 + SEM_2^2}$$

where  $SEM_1$  is the standard error of measurement for the first test and  $SEM_2$  is the standard error of measurement for the second test. Using this formula one person's scores on different tests can be compared. This can be particularly useful when tests are being used to identify an individual's relative strengths and weaknesses, possibly for development purposes. In selection situations it is more common to compare different people's scores on the same test. In this situation,  $SEM_1$  and  $SEM_2$  have the same value, meaning that the formula can be simplified to:

$$SED = 1.414 \times SEM$$

As with the SEM, if two scores differ by one SED or more, the higher scorer is likely to remain on top 68% of the time – about two times out of three. Alternatively, this situation can be expressed as being 68% certain that there is a real difference between the scores. By multiplying the SED by 1.28 or 2.0, the level of certainty can be increased to 80% or 95% that two people's scores really are different. The SED in raw scores and T-scores for each of the Reasoning Tests is shown in **Table 9** below.

Closed tests							
		68% SED		80% SED		95% SED	
		Raw score	T-score	Raw score	T-score	Raw score	T-score
Verbal	1	2.56	4.47	3.28	5.72	5.12	8.94
	2	3.28	6.32	4.19	8.09	6.55	12.65
	3	3.21	5.29	4.11	6.77	6.42	10.58
	4	3.20	5.10	4.09	6.53	6.39	10.20
Numerical	1	1.74	3.74	2.22	4.79	3.47	7.48
	2	2.68	5.66	3.43	7.24	5.36	11.31
	3	2.90	5.10	3.71	6.53	5.80	10.20
	4	3.05	4.69	3.90	6.00	6.10	9.38
Abstract	1	2.93	3.74	3.74	4.79	5.85	7.48
	2	4.20	5.10	5.38	6.53	8.40	10.20
	3	4.47	4.00	5.72	5.12	8.94	8.00
	4	4.42	4.24	5.65	5.43	8.83	8.48
Open tests							
Verbal	1	4.95	4.00	6.33	5.12	9.89	8.00
	2	4.38	4.24	5.60	5.43	8.76	8.48
	C*	2.66	5.66	3.40	7.24	5.32	11.31
Numerical	1	4.30	4.00	5.51	5.12	8.61	8.00
	2	3.55	5.48	4.54	7.01	7.10	10.95
	C	2.20	5.29	2.81	6.77	4.39	10.58
Abstract	1	4.26	3.16	5.45	4.05	8.52	6.32
	2	4.67	4.00	5.97	5.12	9.33	8.00
	C	3.28	5.48	4.20	7.01	6.56	10.95

Table 8: SED for the PfS-Reasoning Tests at 68%, 80% and 95% confidence levels

In order to use **Table 7**, first identify the confidence level required (68%, 80% or 95%) and whether raw scores or T-scores are being used. Find the appropriate column using the first two rows of **Table 7**. Then find the appropriate test in the left-hand column and follow the row across until it intersects with the column to obtain the SED. Test scores need to differ by at least the SED before the difference can be said to be real. For example, to be 80% certain that raw scores from Numerical test Level 2 closed test reflect a real difference in numerical reasoning ability, the difference between raw scores has to be at least 3.44 points. The values in **Table 7** are given as decimals whereas test scores will typically be whole numbers. If users wish to work with whole numbers for simplicity, SEDs should always be rounded up and never down, as rounding down will reduce the confidence that can be placed in the SED. Rounding up will effectively make no difference, as test scores are whole numbers.

## Bias

When used appropriately, psychometric tests have the potential to offer objective, unbiased assessments of ability, aptitude or personal characteristics. Bias in testing occurs because tests have been poorly constructed or because they are used inappropriately. An overview of how to select tests appropriately has been given in **Section Two**. Ensuring that all people are tested under the same conditions by following the standardised administration procedure further reduces the possibility of bias.

More fundamental than appropriate test use is test construction; if a test is inherently biased, the results it gives will always be biased regardless of whether it is being used and administered appropriately. Test bias can arise when the test measures the construct it purports to, but also another, unrelated construct. If the level of this unrelated construct varies between different groups, then the overall results from the test may be biased. For example, a numerical test may contain a large verbal component. If verbal ability differs between two groups (say, people with and without English as their first language) scores may favour one group over another, even if the assessment of numerical ability within the test is fair to both groups.

The initial development of the PfS-Reasoning Tests involved the definition of the areas to be assessed and identification of appropriate test formats (see **Section Four**). From this definition the test specifications were developed, including the descriptions of suitable item content for each Reasoning Test. Bias was therefore minimised by ensuring that the tests did not assess constructs other than the core Verbal, Numerical or Abstract reasoning abilities. Test items were also reviewed for possible bias and subjected to bias analyses during the trailing stage.

Bias can be assessed in two ways: through an examination of overall test scores or the difficulty of individual test items. To assess whether differences in total test scores indicate bias or reflect real differences in the constructs being assessed, it is necessary to find a marker against which test scores can be assessed. As pure markers for constructs such as reasoning abilities are very difficult to identify, the item-level approach to bias was used in the development of the PfS-Reasoning Tests.

The item-level bias analyses conducted during the development of the PfS-Reasoning Tests used a technique known as differential item functioning (*dif*). *Dif* analyses identify whether individual test items are found to be disproportionately easy or hard by different groups of test takers, once their overall score on the test has been allowed for. In other words, if two groups of test takers (say, males and females) obtain very similar overall test scores, the chances of them answering each item correctly should be approximately the same. If one group has a much higher chance of answering an item correctly, the item may be biased.

*Dif* analyses require quite large samples for the results to be robust. Analyses were conducted during the initial stages of development for males and females and for test takers who described their ethnic background as being 'white' and those from other

ethnic backgrounds ('non-whites'). More detailed *dif* analysis of the specific ethnic groups was not possible during the initial stages of development due to the large samples required to do this reliably. Few items were seen to show significant *dif*, suggesting that the pre-trialling reviews and screening of items prior to constructing the final versions of the tests had successfully identified problematic items.

Mean test scores were also examined for males and females and 'whites' and 'non-whites'. The results of these are shown in **Tables 09 and 10** below. As can be seen from **Table 09**, significant score differences were observed between males and females on a number of the PfS-Reasoning Tests. With relatively large sample sizes, however, even very small differences between groups can reach statistical significance. Because of this it is more appropriate to examine differences in terms of 'effect sizes', which look at the difference between groups as a proportion of the pooled standard deviation (taken from **Table 6**). Effect sizes are shown in the last columns of **Tables 09 and 10**.

Guidelines for interpreting effect sizes describe values less than 0.2 as indicating 'small' differences between groups, those between 0.2 and 0.5 as 'medium' and those above 0.5 as 'large' (Cohen, 1988). All differences between males and females fall into the 'small' or 'medium' effect sizes, though there is no consistent pattern of differences between the different test types. This suggests that many of the observed differences may be due to the characteristics of specific samples.

Comparisons of the mean test scores of 'whites' and 'non-whites' revealed a number of statistically significant differences and a number of cases where the effect sizes associated with these differences were of either a medium (10 out of 21 comparisons) or large (5 out of 21 comparisons) magnitude. The remaining 6 comparisons were of a small magnitude. In all cases the differences were seen to favour the 'white' group over the 'non-white' group. These findings reflect the well-established evidence that people from ethnic minority groups, on average, tend to achieve lower scores on ability tests (e.g. CollegeBoard, 2003).

Simple comparisons such as 'whites' and 'non-whites' can mask more complex patterns of performance seen between more precisely defined groups, but to make such comparisons usually requires large numbers of test takers so that all groups are adequately represented. Such comparisons were possible with the open level 2 Reasoning Tests, which have been used extensively in universities and for which large samples of test takers from more precisely specified ethnic groups were available. The results of this analysis can be seen in **Table 11**, which shows the mean test scores according to the 16 ethnic groups defined for the 2001 census of England and Wales.

	Males			Females			Difference	Effect size	
	Mean	SD	Sample size	Mean	SD	Sample size			
<b>Closed Reasoning Tests</b>									
Verbal	1	14.99	5.75	78	17.58	5.51	132	2.59**	0.45
	2	16.20	5.06	158	16.44	5.32	145	0.24	0.05
	3	24.34	6.28	614	23.90	5.88	708	0.44	0.07
	4	25.68	6.46	580	25.22	6.07	551	0.46	0.07
Numerical	1	18.87	5.01	133	19.79	4.16	117	0.92	0.20
	2	15.20	4.57	201	14.57	4.97	136	0.63	0.13
	3	18.69	5.93	864	17.29	5.30	745	1.40***	0.25
	4	16.87	6.71	883	15.35	6.10	627	1.52***	0.23
Abstract	1	27.55	8.15	75	29.41	7.43	81	1.86	0.24
	2	18.82	7.96	125	22.91	8.02	117	4.09***	0.50
	3	30.81	11.66	446	31.61	10.65	414	0.80	0.07
	4	30.64	10.56	503	29.98	10.21	378	0.66	0.06
<b>Open Reasoning Tests</b>									
Verbal	1	17.03	12.21	434	13.29	12.25	576	3.74***	0.30
	2	29.60	10.27	12813	29.62	10.39	11259	0.02	0.01
	C <sup>#</sup>	13.30	4.81	396	14.24	4.45	367	1.06**	0.23
Numerical	1	13.75	11.20	716	15.42	10.18	640	1.67**	0.16
	2	18.67	6.64	20966	17.85	6.23	16275	0.82***	0.13
	C <sup>#</sup>	12.18	4.43	396	12.54	3.83	367	0.36	0.09
Abstract	1	37.85	13.45	301	40.73	13.34	214	2.88*	0.21
	2	33.04	11.99	7349	34.54	11.20	5682	1.50***	0.13
	C <sup>#</sup>	15.66	5.57	396	18.09	6.17	367	2.43***	0.41

#Combined Reasoning Test  
 \*p<0.05; \*\*p<0.01; \*\*\*p<0.001

*Table 09: Mean raw scores and standard deviations for males and females on the PfS-Reasoning Tests*

	'whites'			'non-whites'			Difference	Effect size	
	Mean	SD	Sample size	Mean	SD	Sample size			
Closed Reasoning Tests									
Verbal	1	19.09	6.49	708	16.81	6.27	79	2.28**	0.40
	2	18.80	5.31	233	14.22	6.08	27	4.58***	0.88
	3	25.43	5.60	817	22.25	6.17	228	3.18***	0.52
	4	26.46	5.77	772	23.47	6.94	230	2.99***	0.48
Numerical	1	20.17	5.69	721	19.56	5.42	94	0.61	0.06
	2	15.62	5.67	273	13.72	4.44	36	1.90	0.26
	3	18.79	5.83	989	16.74	5.54	324	2.05***	0.44
	4	16.25	6.35	1028	15.57	6.76	306	0.68	0.10
Abstract	1	28.12	9.80	547	27.97	8.20	61	0.15	0.02
	2	24.18	8.06	155	19.65	5.51	20	4.53*	0.55
	3	31.80	11.25	531	28.96	10.95	102	2.84***	0.25
	4	30.60	11.00	641	26.78	10.02	141	3.82***	0.37
Open Reasoning Tests									
Verbal	1	22.19	11.98	388	10.05	10.21	564	12.14***	0.98
	2	32.22	9.46	14389	25.54	10.08	8622	6.68***	0.65
	C#	13.94	4.65	680	12.05	4.91	76	1.89***	0.42
Numerical	1	16.27	10.65	675	12.82	10.60	681	3.45***	0.32
	2	18.37	6.06	22349	18.20	7.02	13372	0.17*	0.03
	C#	12.47	4.10	680	11.39	4.51	76	1.08*	0.26
Abstract	1	41.05	12.54	319	35.67	13.78	180	5.38***	0.40
	2	34.47	11.03	8124	32.26	12.34	4346	2.21***	0.19
	C#	16.98	5.84	680	15.87	7.05	76	1.11	0.19

#Combined Reasoning Test  
 \*p<0.05; \*\*p<0.01; \*\*\*p<0.001

*Table 10: Mean raw scores and standard deviations for 'whites' and 'non-whites' on the Pfs-Reasoning Tests*

	Verbal				Numerical				Abstract			
	N	Mean	SD	Effect size	N	Mean	SD	Effect size	N	Mean	SD	Effect size
British	10499	32.83	9.39		15456	18.46	5.96		5412	34.88	10.76	
Irish	1213	31.84	9.01	0.10	2029	17.78	5.44	0.10	641	33.30	10.80	0.14
Any other White background	2677	30.00	9.60	0.27	4864	18.34	6.57	0.02	2071	33.78	11.72	0.09
White and Black Caribbean	144	19.83	11.84	1.26	199	13.54	7.07	0.76	125	27.59	11.21	0.62
White and Black African	62	29.13	11.12	0.36	108	15.59	5.67	0.44	61	29.66	13.42	0.45
White and Asian	247	28.38	12.13	0.43	321	18.80	6.48	0.05	120	36.78	13.00	0.16
Any other mixed background	246	29.99	10.39	0.28	391	17.96	6.30	0.08	150	34.10	10.64	0.07
Indian	2126	26.82	9.93	0.58	3715	17.25	6.67	0.19	1132	30.14	12.19	0.41
Pakistani	406	27.44	8.88	0.52	643	16.68	6.18	0.27	168	29.97	10.70	0.42
Bangladeshi	184	25.86	9.47	0.68	262	17.29	6.60	0.18	76	33.67	10.73	0.10
Any other Asian background	760	25.79	9.89	0.68	1154	17.75	6.99	0.11	417	31.38	10.45	0.30
Caribbean	198	24.40	9.14	0.82	308	14.80	6.31	0.56	98	28.48	10.49	0.55
African	656	24.12	9.17	0.84	1090	15.84	6.00	0.40	360	27.25	9.76	0.65
Any other Black background	65	25.58	7.52	0.70	90	15.40	5.98	0.47	31	26.84	9.18	0.69
Chinese	3289	24.27	9.97	0.83	4718	20.38	7.16	0.30	1463	35.59	12.96	0.06
Any other	239	27.20	10.65	0.55	373	18.07	7.12	0.06	145	35.17	12.49	0.02

*Table 11: Mean test scores and effect sizes for different ethnic groups based on the open Level 2 PFS-Reasoning Tests Reasoning Tests*

The detailed analysis of test scores obtained by different ethnic groups shown in **Table 11**, indicates the means and SDs for each group on the Verbal, Numerical and Abstract level 2 open tests. The effect size for each group is also shown, and indicates the extent to which the mean for each group differs from the 'British' mean. The 'British' group obtained the highest mean score on the Verbal test, with the 'Chinese' group obtaining the highest mean score on the Numerical and Abstract

tests. In terms of lower scoring groups, 'White and Black Caribbean', 'Caribbean', 'African' and 'Any other Black background' consistently showed some of the largest effect sizes.

With the recent introduction of legislation on age, there has been particular interest in how performance on tests of mental ability such as the PfS-Reasoning Tests is influenced by age. The links between test scores and age are seen in **Table 12**. The strongest links between age and test performance are seen the closed level 1 tests and the Combined test. These positive correlations indicate that test scores increase as does respondents' age, and most likely reflects progress and development through the education system as these tests were taken by people as young as 14. The largest of these associations, seen for the closed level 1 Verbal test, indicates that age accounts for just over 8% of the variance in test scores. Some evidence of a negative association with age was also seen amongst samples with slightly older respondents, though there was no evidence of a highly significant fall-off in performance. Taken across all test levels, age accounted for less than 2% in performance on average.

Test version	Verbal	Numerical	Abstract
Closed Tests			
1	0.29 (n=803)	0.23 (n=831)	0.04 (n=598)
2	0.03 (n=549)	0.15 (n=595)	-0.15 (n=429)
3	-0.06 (n=1340)	-0.10 (n=1624)	-0.12 (n=878)
4	0.13 (n=1137)	0.05 (n=1552)	0.02 (n=889)
Open tests			
1	0.22 (n=1010)	-0.10 (n=1343)	-0.13 (n=512)
2	-0.02 (n=23999)	-0.13 (n=37134)	-0.13 (n=12981)
Combined	0.22 (n=755)	0.18 (n=755)	0.14 (n=755)

*Table 12: Associations between raw PfS-Reasoning Tests and respondents age*

## A commentary on interpreting bias data

When variations in test scores are seen between different groups, whether those groups are defined on the basis of sex, ethnicity, age or any other factor, an immediate possibility is that the tests are biased. That is, in some way, the items within the test or the whole testing process itself are easier for some groups than others, resulting in differential performance. Such differential performance is not in itself a bad thing, but becomes an issue if it can be shown that the differences arise due to test performance being affected by extraneous factors unrelated to the construct being assessed – in the current case, Verbal, Numerical or Abstract reasoning ability.

With a focus on the analysis of ethnic groups performance, as it is here that the largest mean differences were seen, the purpose of this section is to explore the possible reasons for this differential test performance. It is recognised, and largely accepted, that variations in test performance will be seen between different ethnic groups. These differences remain despite the best efforts of test developers to make tests fair and accessible through careful test design, item writing and review, trialling and item-level statistical analyses. Comparisons between the current tests and other aptitude tests will therefore indicate the extent to which the PfS-Reasoning Tests can be considered as functioning within ‘accepted’ parameters.

The Scholastic Assessment Test (SAT) used as part of college selection in America and other countries for over 2 million students each year, is probably the most researched and well-developed aptitude testing programme, and so provides an appropriate benchmark for the examination of ethnicity. ETS, who develop the SAT, have also been influential in shaping modern thinking on test bias and how to identify it.

In terms of effect sizes, a difference of 0.98 (almost 1 SD) is seen between ‘White’ and ‘African American’ candidates on the verbal part of the SAT and a difference of 1.08 seen on the math part, with ‘African American’ candidates scoring lower in both cases. When compared to the ‘Asian American’ group, ‘White’ candidates score 0.21 higher on verbal, but 0.41 lower on math (CollegeBoard, 2003). These figures indicate that substantial differences between the mean scores of different ethnic groups remain, despite the best efforts of test developers. They are also in line with the findings from **Table 11**, where the ‘White and Black Caribbean’, ‘Caribbean’, ‘African’ and ‘Any other Black background’ groups showed some of the largest effect sizes when compared with ‘Whites’ and obtained some of the lowest mean test scores.

It is currently unclear why these differences are seen, although there are a number of possibilities (see for example Freedle, 2003 and Neisser, Boodoo, Bouchard, Boykin, Brody, Ceci, Halpern, Loehlin, Perloff, Sternberg and Urbina, 1996 for a discussion). First, any differences may reflect true differences in the capability of candidates. As there is no ‘gold standard’ against which aptitudes can be measured, it is very difficult to establish any individual’s or group’s true level of specific abilities. The only robust way of checking a test for bias, and so determining whether test score differences

reflect differences in the ability to perform a job, is through comprehensive validity studies. Second, the possibility of 'differential sampling' needs to be considered. The effect of any background factors on test scores could be due to groups being made up from people of different ability levels. To determine whether differential sampling is affecting the observed scores it would be necessary to collect additional background information on candidates, particularly educational qualifications and proficiency in English. Third, differences could be due to variations in familiarity with reasoning tests. Finally, a number of cultural and sociological arguments have been proposed to explain differential test performance (see Neisser *et al* 1996 for a summary). Many of these theories focus on the meaning and experience of testing to people from different cultures, recognising that the whole testing movement has its roots in a 'white, middle-class' philosophy.

To summarise, differences in the mean test scores of different groups do not prove that a test is biased. The ethnic differences observed in the PfS-Reasoning Tests are also seen in other widely used tests and remain despite intensive efforts to make tests 'fair'. Ensuring the fairness of the tests is an ongoing project combining test research, support to candidates and the need to validate the tests against meaningful and reliable job-related criteria. Further analyses will therefore be conducted and reported in subsequent versions of this User's Guide when sufficient data is available.

## Validity

Validity is the most important consideration when using any test or assessment. If a test is valid, it will produce meaningful results and will contribute significantly to the decision-making process, either predicting subsequent job or training performance or correctly identifying development needs.

Many different forms of validity have been identified, but it is most accurately viewed as a unified concept (Messick, 1995), with different forms of validity contributing to an overall judgement. Further, it can never be asserted that a test is globally valid or not, as validity relates to the use of a test in specific situations – its 'fitness for purpose'. Four main types of validity are discussed here: face, content, construct and criterion validity.

### Face validity

A test has face validity when it looks as though it is measuring what it claims to measure. Although not always considered to be a genuine source of validity, if test takers can clearly see links between the skills being measured and a certain job, they are likely to be motivated to complete the test to the best of their abilities. There may be lower motivation to perform well if the reasons for completing the test are unclear. Further, the selection process is seen as a form of social interaction, during which applicants will form impressions of an organisation (Anderson and Cunningham-Snell, 2000). The use of tests with low face validity may have a negative impact on this emerging impression. For these reasons, face validity is important when using assessments in occupational settings.

Evidence for the face validity of the PfS-Reasoning Tests was collected during the trialling stage by observing test sessions and obtaining feedback from test takers. Users found the tests easy to use and the content to be acceptable. Further, the feedback reports designed for test takers were also seen to be accessible, informative and to provide useful points for consideration. However, some users found the reports to be quite long and suggested that reports simply containing the test results would have been more useful for their purposes. To address this need, summary reports for both administrators and test takers were created (see **Appendix Two** for a sample). Although feedback indicated that the tests had good face validity, this has to be supported by other forms of validity or a test may be accused of superficiality.

### **Content validity**

If the items in a test provide adequate coverage of the area being assessed, and do not relate to areas outside the sphere of the test, then the test is said to have content validity. For the Verbal, Numerical and Abstract Reasoning Tests, the process of ensuring content validity started by developing the test specifications, detailing the content of each test. The review of test items by the test development team and external experts further contributed to content validity, ensuring that items met the test specifications and making necessary changes where they did not. The final stage in this process was the compilation of the tests themselves, where the content of each separate test was carefully checked to make sure it was adequate.

The development process has resulted in tests which fulfil the specifications set out in **Section Four**. Ultimately, however, potential test users should review the tests themselves, to ensure that test content sufficiently matches their needs.

### **Construct validity**

Construct validity refers to what a test actually measures. In the case of the PfS-Reasoning Tests, the constructs are Verbal, Numerical and Abstract reasoning ability. Evidence for construct validity comes from the examination of how scores on each of the tests relate to each other and to established assessments that measure related constructs.

The correlations between the three PfS-Reasoning Tests at each level are shown in **Table 13**. A number of observations can be made from this data. Firstly, the correlations show that each of the Reasoning Tests is assessing a quite distinct area of ability. The highest correlation is between the Numerical and Verbal parts of the Combined Test, showing that the two share around 42% of common variance (i.e. performance on one test will account for no more than 42% of the performance on another). Among the closed tests the degree of association is generally far less, with the mean correlation indicating that just under 20% of common variance is shared between tests.

Secondly, there is a decrease in the mean correlations between the higher levels of the closed tests. The mean correlations are 0.56, 0.45, 0.45 and 0.32 for levels 1 to

4 respectively. It is known that as people get older and specialise in their areas of study abilities tend to become more defined, meaning that the correlations between assessments of different abilities are reduced. The pattern of relationships found with the PfS-Reasoning Tests supports this differentiation of abilities. This observation is further supported by the data from the Combined Reasoning Test, where the majority of test takers were in the last two years of compulsory education. The level of correlation in this test is also likely to be influenced by the fact that the three sub-sections are relatively short and taken immediately after each other, so potentially reducing some of the sources of error in test scores.

Together these findings show a meaningful pattern of relationships within the three PfS-Reasoning Tests, indicating that they assess quite distinct areas of reasoning ability and so supporting the validity of the constructs defined for the PfS-Reasoning Tests.

### Closed Reasoning Tests

	Verbal 1	Numerical 1
Numerical 1	0.60 (189)*	-----
Abstract 1	0.53 (131)	0.54 (120)

	Verbal 2	Numerical 2
Numerical 2	0.53 (263)	-----
Abstract 2	0.41 (237)	0.39 (241)

	Verbal 3	Numerical 3
Numerical 3	0.48 (1659)	-----
Abstract 3	0.48 (1304)	0.40 (1218)

	Verbal 4	Numerical 4
Numerical 4	0.28 (1240)	-----
Abstract 4	0.35 (805)	0.33 (813)

### Open Reasoning Tests

	Verbal 1	Numerical 1
Numerical 1	0.37 (1288)	-----
Abstract 1	0.38 (106)	0.44 (118)

	Verbal 2	Numerical 2
Numerical 2	0.38 (5820)	-----
Abstract 2	0.47 (3565)	0.35 (3731)

	Combined - Verbal	Combined - Numerical
Combined - Numerical	0.65 (880)*	-----
Combined - Abstract	0.56 (880)	0.62 (880)

\* figures in parentheses indicate number of test takers

*Table 13: Intercorrelations of the PfS-Reasoning Tests*

Further evidence for the criterion validity of the PfS-Reasoning Tests comes from a number of studies that have explored the association between them and other assessments of capability. These studies are summarised below.

Association between level 4 closed tests and the Graduate Management Admission Test (GMAT), which is used by graduate management schools in many countries as part of their admission process, was examined in a sample of postgraduate students at a business school based in London during 2004. The sample consisted of approximately 56% males and 44% females, with a mean age of 26.39 years

(SD=4.47). A significant proportion of the students in this sample came from outside of the UK, though exact data on this was not available.

As shown in **Table 14**, the strongest association was seen between GMAT and the PfS Verbal test. This would be expected as the GMAT contains two sections of verbal material and one of numerical. There is no equivalent in the GMAT to the PfS Abstract test, as reflected in the lower association of PfS Abstract with the GMAT scores. It should also be noted that respondents in this sample were asked to recall their GMAT scores from memory and that the time between taking the two assessments could have been around a year for some students, both of which could have affected the resulting correlations.

	Correlations with GMAT and sample size
Verbal 4	0.34 (n=74)
Numerical 4	0.23 (n=97)
Abstract 4	0.15 (n=54)

*Table 14: Associations between PfS-Reasoning Tests and the GMAT*

The Graduate and Managerial Assessment (GMA; Blinkhorn, 1985) is an established and widely used test, consisting of high-level verbal, numerical and abstract tests. As with the PfS Abstract Reasoning Tests, the GMA Abstract test is based on Bongard problems (see page 32), and this study explored the association between levels 3 and 4 of the closed PfS Abstract Reasoning Tests and the GMA Abstract Test form A (GMA-A). Data was collected during the first quarter of 2007 from two groups of Year 12 students, one at a boys-only comprehensive school and another at a girls-only independent school. There were 78 participants from the boys school with a mean age of 16.7 years (SD=0.7) and 48 from the girls school with a mean age of 16.4 (SD=0.5). The order of test completion was counterbalanced.

The correlations between the scores from the three tests are shown in **Table 15**, with the first figure showing the raw correlation and the second in brackets the correlation corrected for the reliability of the two tests in question. All uncorrected correlations between PfS tests and GMA are 0.64 or greater, and when corrected for reliability are 0.71 or greater. These figures indicate a good degree of association between the assessments and all exceed the 0.70 threshold, typically recognised the point at which tests can be considered to be alternate forms of each other.

	PfS 3	PfS 4	GMA lenient
PfS 4	0.73 (0.80)		
GMA lenient	0.71 (0.78)	0.64 (0.71)	
GMA harsh*	0.68 (0.86)	0.69 (0.79)	0.80 (0.89)

\* Harsh scoring on the GMA-A awards one mark only if all of the five test shapes in a group have been answered correctly.

*Table 15: Associations between PfS Abstract Tests and GMA Abstract form A*

The association between the PfS Verbal and Numerical level 1 open tests and SHL's VMG3 (verbal reasoning) and NMG3 (numerical reasoning) was examined in a sample of employees at a UK emergency services organisation. Forty-four employees completed both verbal tests (mean age 42.32, SD=5.02) and seventy the numerical tests (mean age 42.26, SD=5.31). All test takers were male.

The employees completed the PfS tests as preparation for an internal development programme and the SHL tests subsequently as part of the programme. The correlations between the verbal tests were 0.43 (0.51 corrected for reliability) and between the numerical tests were 0.26 (0.29 corrected for reliability).

Further data relating to construct validity was obtained from versions of the PfS-Reasoning Tests constructed with items from each of the four closed test levels. These tests were developed for a client who needed to assess people across a wide range of ages and ability levels. The correlations between the versions of the PfS-Reasoning Tests and the client's existing assessments were examined for evidence of construct validity. As the reliability of the client's existing assessments was poor the correlations in **Table 16** have all been corrected for reliability.

The data reported in **Table 16** was collected in 2003 from 254 candidates to the client's organisation. The mean age of the candidates was 22.2 years (SD=3.0), and 240 (94.5%) were male and 14 (5.5%) were female.

	Correlations with existing reasoning tests and sample size
Verbal reasoning correlations	0.48 (121)
Numerical reasoning correlations	0.65 (115)
Abstract reasoning correlations	0.36 (122)

*Table 16: Intercorrelations between the Verbal, Numerical and Abstract Reasoning Tests and existing reasoning tests*





































































































