

# Inter-rater reliability of the Norwegian version of the Perceive, Recall, Plan and Perform System of Task Analysis

AV LINDA STIGEN AND JULIE PAGE

## Abstract

### *Background/aim*

In order to survive in the current health and social care climate, occupational therapists need to demonstrate that their interventions are effective both clinically and economically. The use of the Perceive, Recall, Plan and Perform system of task analysis is expanding in Norway and there is a need to be sure that we are using this assessment in a valid and reliable way. The aim of this study was to establish inter-rater reliability of the Norwegian version of the Perceive, Recall, Plan and Perform System of Task Analysis.

### *Methods*

Inter-rater reliability was measured using percentage agreement (PA) and intra-class correlation coefficient (ICC).

### *Results*

29 occupational therapists participated and rated three clients. ICC results ranged between 0.23 to 0.98 on sub-quadrant score and 0.15 to 0.62 on quadrant score. The PA ranged from 42 percent to 100 percent on descriptor scores, and between 77 percent and 83 percent on quadrant level.

### *Conclusions*

The results indicate that the inter-rater reliability of this group of occupational therapists is low compared to previous studies. The study gives some indication on the inter-rater reliability of the Norwegian version, however future studies are necessary in order to be able to draw further conclusions.

*Keywords;* psychometrics, reliability of results, occupational therapy, cognition, task performance and analysis, (PRPP)



*Linda Stigen, MSc, is an OT specialist in somatic health and works with community based rehabilitation in Bjerke district in Oslo. The article is based on her master thesis in the European Master of Science in Occupational Therapy program. The project was done while she was working at and with support from Sunnaas Rehabilitation Hospital. Email: lindastigen@hotmail.com*



*Julie Page, PhD, is head of Research and Development, Institute of Occupational Therapy, School of Health Professions, Zurich University of Applied Sciences, Switzerland*

*Denne artikkelen ble mottatt 19. februar 2011  
Den ble akseptert for publisering 6. februar 2012  
Oppgitte interessekonflikter: Ingen*

## Introduction

Impairments of cognitive function are a significant cause of disability after traumatic brain injury (TBI) and stroke (Cicerone et al, 2000). Within the context of occupation, cognitive impairments are likely to have impact on some aspects of life, and occupational therapy is an important part of rehabilitation of these impairments (Grieve & Gnanasekaran, 2008). In evidence-based occupational therapy one of the most important underpinnings is the consistent use of outcome measures to evaluate occupational therapy services (Law, Baum & Dunn, 2005). Understanding the psychometric properties of standardised assessments is a fundamental skill required by occupational therapists today (Laver Fawcett, 2008). In order to survive in the current health and social care climate, occupational therapists need to demonstrate that their interventions are effective both clinically and economically (Taylor, 2007; Laver Fawcett, 2008). Measurements in occupational therapy are used to improve decisions regarding specific clients or programs, and information gained through measurement helps occupational therapists design interventions for individuals and evaluate the outcome of these programs (Law, Baum & Dunn, 2005). Assessments of cognitive function can be considered to be the starting point of occupational therapy interventions (Grieve & Gnanasekaran, 2008) and therefore one needs appropriate, valid and reliable outcome measurements (DePoy & Gitlin, 2005). However, there are currently only a few valid and reliable instruments existing in other languages than English. This means that clinicians working in non-English speaking countries can choose between not using standardised instruments, creating new instruments, or translating instruments from English into their native language. The first option is a challenge, as clinicians are not able to refer to normative data with confidence when assessing clients. The second option is often not possible because of time constraints and the expensive financial and technical resources. Occupational therapists are therefore often left with the option to translate previously developed instruments into their own language. However, when international or cross-cultural applications are involved in therapy or in research, one of the major problems that clinicians and researchers encounter is the translation of material (Beaton, Bombardier, Guillemin & Ferraz, 2000). Historically, occupational therapists have used non-standardised assessments, especially unstructured interviews and observations, or standardised assessments were modified to suit their different clinical environment (Laver Fawcett, 2008). This is due to language issues but also a number of other reasons; a lack of appropriate standardised assessments, poor resources; that standardised assessments can be lengthy to administer; that non-standardised assessments are flexible in terms of procedures and settings, etc. A common problem has also been that of taking different parts of standardised tests or individual test items to integrate these into a «therapist-constructed» assessment battery. However, once the standard procedure for test administration and scoring has been changed, there is no longer any guarantee of the reliability and vali-

dity of the test (Turner, Foster & Johnson, 2002).

## The Perceive, Recall, Plan and Perform System of Task Analysis

The Perceive, Recall, Plan and Perform System of Task Analysis (PRPP) is a two stage, standardised, criterion-referenced assessment (Chapparo & Ranka, 2006). A criterion-referenced assessment means that the client is measured up against predefined criteria instead of a normative sample. This is useful when a person's competence or level of mastery is examined (Turner, Foster & Johnson, 2002). PRPP is an occupation-focused, client-centred assessment of occupational performance. The PRPP evaluates information processing in daily activities as well as level of mastery (Chapparo & Ranka, 2006). Stage one focuses on the assessment of mastery in occupational performance, and it gives an idea on how well the person performs a task in addition to giving us the type of errors the person makes. The task is broken down into steps where four error types are identified in each step. The error types are errors of accuracy, errors of repetition, errors of omission and errors of timing (Chapparo & Ranka, 2006). Stage two is used to identify the information processing that interferes with performance. The information strategies identified in the PRPP include attention and sensory perception (Perceive), memory (Recall), response planning and evaluation (Plan) and performance monitoring (Perform). The information processing strategies are divided into 34 error types which are called descriptors. The descriptors are the items that are identified as the clients perform the tasks and are recorded in stage two of the assessment. The stage two model is presented in figure 1.

The PRPP System of Task Analysis has proven to be a good tool to assess level of function of people with traumatic brain injuries (Fry & O'Brian, 2002). It has also proven to be a good tool in order to measure change in cognitive abilities in activities of daily living (Nott & Chapparo, 2007). Two studies investigating inter-rater reliability of the PRPP have previously been published. One was done in Canada by Aubin et al in 2008, and one was done in Australia by Nott, Chapparo & Heard in 2008 (Aubin et al, 2008; Nott, Chapparo & Heard, 2008). The results of these studies are presented in table 1.

The use of PRPP is expanding in Norway, and therefore there was a need for a Norwegian translation in order to ensure the validity and reliability when using the PRPP in the Norwegian context.

## Measurement of inter-rater reliability

Reliability refers to the consistency of measurements when the testing procedure is repeated on a population of individuals or groups (American Educational Research Association APA, 1999). Reliability is an interaction among the instrument, the specific group of people taking the test, and the situation (Streiner & Norman 2008). Reliability studies of assessment tools are necessary in order to be sure that the error involved in measurement is small enough to detect changes in what is being measured (Rankin & Stokes, 1998). Whether a scale that has been shown to be

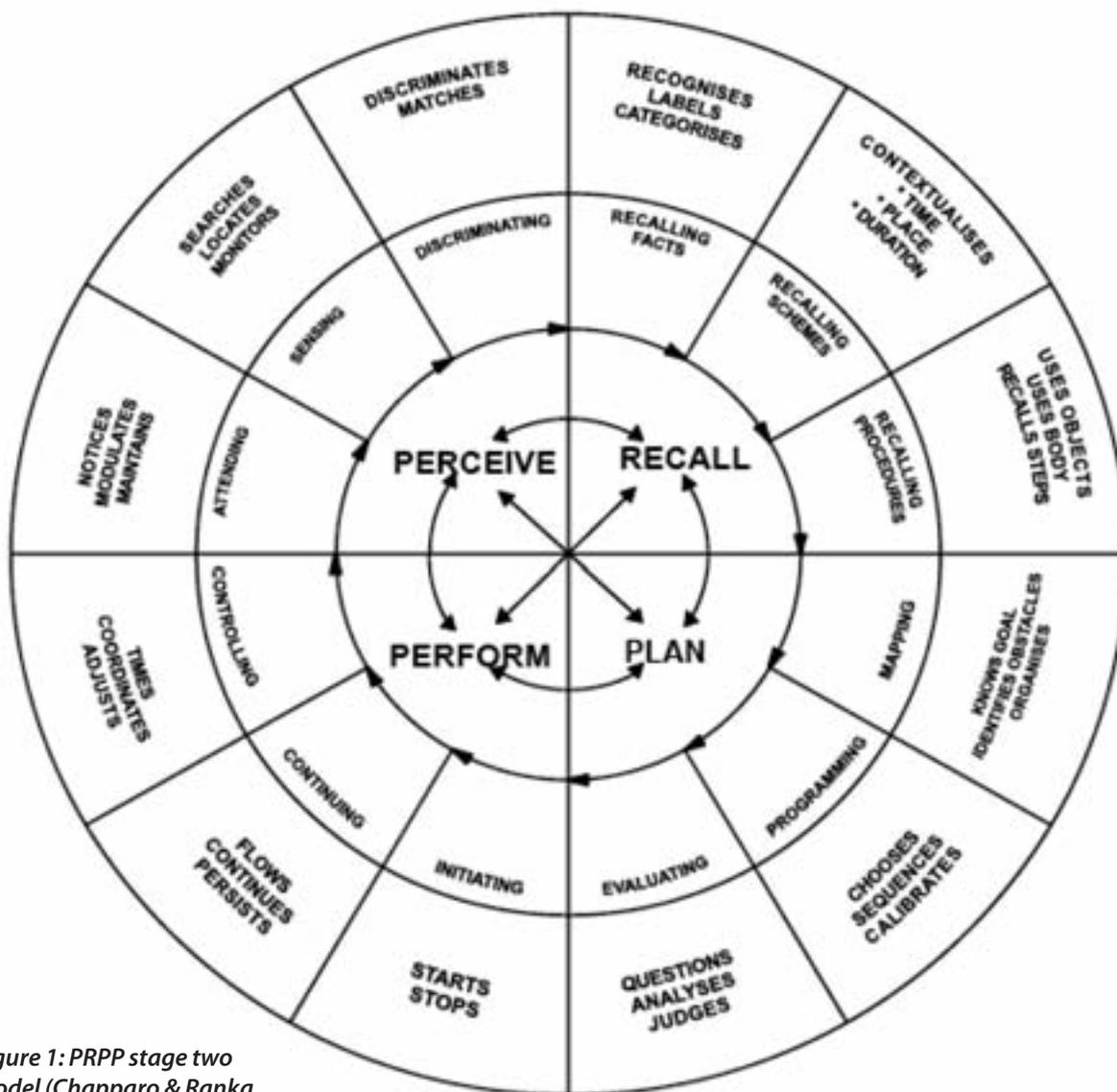


Figure 1: PRPP stage two model (Chapparo & Ranka, 2006).

reliable with one group of individuals in a particular context is reliable with other people and in other situations, must be assessed (Vacha-Haase, 1998; Thompson 2003; Streiner & Norman 2008). It has been said that reliability refers to the results obtained with an evaluation instrument, not to the instrument itself (Gronlund & Linn, 1990; Thompson 2003; Streiner & Norman 2008) so therefore it is more accurate to speak of reliability of test scores rather than of tests (Vacha-Haase, 1998; Thompson 2003; Streiner and Norman 2008). As the degree of established reliability informs the therapist how accurately the scores obtained from an assessment reflect the true performance of the test, it is important for occupational therapists to consider the reliability of an assessment (Turner, Foster & Johnson, 2002).

One type of reliability measurements examines the reliability or consistency of the rater or raters of the test. Inter-rater reliability examines whether different raters can assign similar ratings to the client performing the task. To measure inter-rater reliability the scores of two or more raters are being compared (Laver Fawcett, 2008). Measuring inter-rater reliability is important as many cli-

ents move between occupational therapy services, e.g. from an acute to a rehabilitation ward. As a result one client might be given the same assessment from different therapists (Turner, Foster & Johnson, 2002).

It is important for occupational therapists to consider the reliability of an assessment, especially when it comes to evaluation of the effectiveness of a treatment programme. It is necessary to be sure that we are measuring the true performance and that potential changes in a person's performance on a test are not affected by the time or by the rater (Laver Fawcett, 2008).

The aim of this study was to investigate inter-rater reliability of the Norwegian version of the PRPP.

## Method

### Participants

To investigate the inter-rater reliability of the Norwegian version of the PRPP, 44 occupational therapists were approached during a PRPP training course with the question of participating in the study. Inclusion criteria were that they were present and able to fill in at least one video case, and

	Reliability scores of the different quadrants				Total
	Perceive	Recall	Plan	Perform	
Aubin et al. 2008	0.65	0.65	0.69	0.63	0.77
Nott, Chapparo & Heard, 2008	0.59	0.59	0.51	0.53	0.60

Table 1: Previous studies reporting inter-rater reliability of the PRPP.

	Rater 1	Rater 2	Rater 3	Rater 4	Rater 5	Rater 6	Rater 7	Rater 8	etc
Client 1	48	53	39	51	52	44	44	45	
Client 2	50	48	40	53	49	46	50	51	

Table 2: Example of ICC calculation on total score.

that their mother tongue was Norwegian. 29 occupational therapists agreed to participate and signed an informed consent form. 24 participants filled in three cases and two participants filled in two cases. This gives a dropout percentage of 10 percent. They voluntarily shared that the reasons for this varied: illness, fear of not performing well enough. The video cases used in this study were presented to the participants on day four and five in the training course. The clients were Australian and the videos are normally used in the PRPP training courses. Of the clients in the video cases, two suffered from traumatic brain injury and one client from stroke. One of the clients with traumatic brain injury also had a diagnosis of depression. The activities performed by the clients varied. The task of client one was to put on a sweater, client two was making a cup of coffee and client three was buttering a piece of bread. Each client was assessed on only one task. Video cases were chosen as the method as there is a chance that with live clients, the true performance could be influenced by other issues like becoming nervous when performing the task in front of 26 occupational therapists watching them.

The descriptor *Matches* was not scored in either of the cases as the tasks did not demand it. The mean age of the participants was 37 years (standard deviation (SD) = 10), range 24 - 54. The gender distribution was four percent male (n=1) and 96 percent female (n=25). The mean working experience of the participants was twelve years (SD= 8) with a range of one to 28 years. 96 percent of the participants (n=25) worked with adult neurology, mainly with stroke and TBI clients, while four percent (n=1) worked in education. All the participants attended the same PRPP training course and received the same information on how the PRPP should be applied to clients.

**Data analysis**

Percentage agreement and Intra-class Correlation Coef-

ficient (ICC) were used to analyse the data. Percentage agreement was used as it is the simplest measure of consistency and often applied to criterion-referenced tests that use ordinal scales involving mastery decisions (Laver Fawcett, 2008). ICC was chosen as it can account for multiple sources of error at the same time, and it is easier to use when the group of raters is larger than two (Streiner & Norman, 2008). In addition the ICC reflects the extent of agreement between raters, including a systematic disagreement which is not reflected in other estimates of reliability (Ottenbacher & Tomchek, 1993).

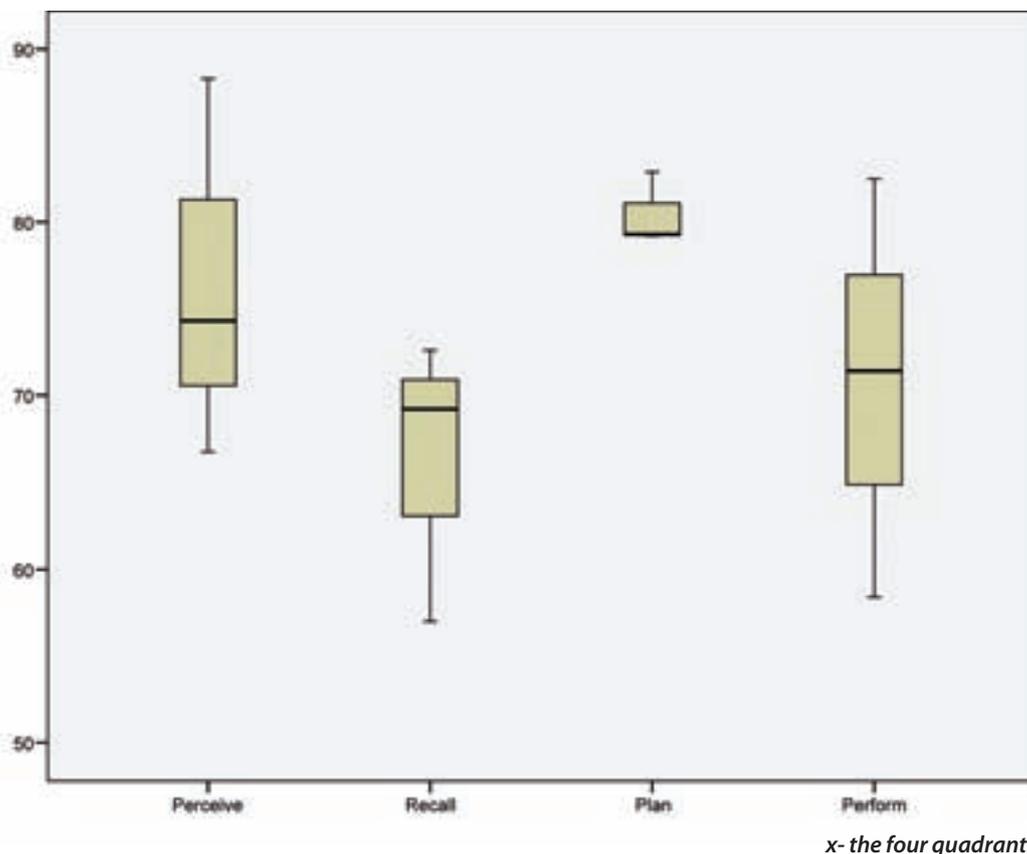
The ICC calculates the variability between the rows and columns, which in this study are the raters and the clients. If the variation within one row is bigger than between columns, the ICC will be low. If the variation within one row is smaller than between rows, the ICC will be high (Streiner & Norman, 2008). The ICC for calculating inter-rater reliability was concerned with the variation between the raters who measured the same group of people. Both raters and clients were considered random, and a two-way random model with absolute agreement was used to obtain the ICC (McGraw & Wong, 1996). SPSS was used for the two-way ICC analysis and the data were set up with raters in the columns and clients in the rows (McGraw & Wong, 1996). The sums of subquadrant, quadrant and total score were used to calculate ICC. An example of the data setup in SPSS is given in table 2.

**Results**

**Percentage agreement (PA)**

PA was calculated on individual descriptor level on the three cases looking for absolute agreement. Based on the percentage agreement on the descriptors, calculations to find the mean on quadrant level was done, and this is presented in Figure 2. The quadrant level ranges from 66.3

y- agreement in percentage



**Figure 2: Percentage agreement on the four quadrants.**

percent on Recall to 80 percent on the Plan quadrant.

The percentage agreement on the descriptor scores ranges from 42 percent on *Categorizes* to 100 percent on *Modulates, Searches and Locates*. 42 percent indicates that 42 percent of the OT's agree on the scoring of that descriptor, and 100 percent indicate that everybody agreed on that descriptor. The mean on the three cases are ranging from 54 percent on *Recognizes* to 97.4 percent on *Locates*. Details on PA for the descriptors are shown in descending order in Figure 3. There is a big range in the PA on descriptor score and there does not seem to be a pattern regarding low PA on descriptors in certain quadrants.

### Intra-class correlation coefficient (ICC)

The ICC for inter-rater reliability was calculated using a two-way random model where the raters and the clients were considered random. ICC on subquadrant levels ranged from 0.23 on *Recalling schemes* to 0.94 on *Controlling*. The quadrant level ranges from 0.15 on *Perceive* to 0.62 on *Recall*. Table 3 shows ICC on subquadrant scores, quadrant scores and on the total score.

The ICC on quadrant level is lower than reported in previous studies. There does not seem to be a pattern regarding whether certain quadrants have subquadrants scores that are lower than others.

## Discussion

### Discussion related to obtained results

This is the first study investigating the inter-rater reliability of the Norwegian version of the PRPP, using data from

stage two. It is important to be aware that the frame of reference when judging reliability scores varies among different authors and between different forms of reliability. It is said that an inter-rater reliability of 0.90 or above is considered to be acceptable level, both for PA and ICC (Benson and Clark, 1982)

The PA in this study was calculated on absolute agreement and has a range from 0.41 (41 percent) on *Categorizes* to 1 (100 percent) on *Modulates, Searches and Locates*. The sample in this study was bigger than in previous ones investigating the PRPP, and thus the results are less likely to arise due to chance (Birkimer & Brown, 1979). However, an agreement rate of 41 percent is considerable lower than what one wants it to be.

The ICC was calculated on subquadrant, quadrant and total score level. This is the first study to report ICC on subquadrant scores. By doing this the results will give information on whether there are specific subquadrants that give lower ICC than others.

### Comparison with previous PRPP inter-rater reliability studies

Several studies that have been published investigated inter-rater reliability of the English version of the PRPP (Aubin et al, 2008; Nott, Chapparo & Heard, 2008). The two studies were performed in Australia and Canada and they both report moderate to good inter-rater reliability based on ICC calculations on total and quadrant scores. The model used for analysing the data varied between these studies; the Australian study used a three-way model where the three factors were raters, clients and tasks (Nott, Chapparo & Heard, 2008) and the Canadian study used a



<b>Perceive</b>	<b>0.15</b>
<i>Attending</i>	0.49
<i>Sensing</i>	0.41
<i>Discriminating</i>	0.94
<b>Recall</b>	<b>0.62</b>
<i>Recalling facts</i>	0.82
<i>Recalling schemes</i>	0.23
<i>Recalling procedures</i>	0.87
<b>Plan</b>	<b>0.26</b>
<i>Mapping</i>	0.75
<i>Programming</i>	0.84
<i>Evaluating</i>	0.69
<b>Perform</b>	<b>0.29</b>
<i>Initiating</i>	0.81
<i>Continuing</i>	0.94
<i>Controlling</i>	0.98
<b>Total</b>	<b>0.67</b>

**Table 3: Descriptive statistics for the intra-class correlation coefficient on subquadrant, quadrant and total score.**

as an issue in the Canadian study (Aubin et al, 2008).

To improve rating differences in observational assessments it has been suggested that there is a need for more extensive training of the raters, and precise clarifications of definitions in addition to scales that are more refined (Slagle et. al, 2002).

### Strengths of the research

This research is the first to investigate properties of the Norwegian version of the PRPP and adding to the international body of knowledge concerning the PRPP. As occupational therapists we are evaluating the occupational performance of our clients, and we need to know that what we are measuring reflects the actual performance.

### Recommendations for future development and research on the PRPP

The results of this study indicate that the inter-rater reliability of the PRPP in this group of occupational therapists is low on quadrant scores compared to previous studies. In order to improve this, it could be important to conduct training in order for the OT's to become more confident in the differences of the descriptors, in addition to conducting workshops to let the occupational therapists watch videos and discuss scoring of clients. It could also be important to work on the criteria for scoring on certain descriptors in order to make them more clear and understandable, making the occupational therapists more capable of differentiating them. As regards whether the Norwegian version should be used by the therapists in Norway or not, it would be important to be well acquainted and used to the English version first. It might be a good idea to have a one-day introduction to the Norwegian version before it is used, in order to get insight into how the version was developed and how it differs from the English version. It could also be a good idea for

the whole training course to be taught in Norwegian, however at this point that is not possible as there are no PRPP trainers in Norway. This is the first study done on the PRPP in Norway and as it is to be used in the Norwegian context with Norwegian clients, it is of high importance to conduct more studies of the psychometric properties of the Norwegian version of the PRPP. Further studies on reliability and also on addressing validity issues related to a wide range of clients should be undertaken.

### Conclusion

This study adds to the growing body of research investigating the psychometric properties of the PRPP. A few studies have been published concerning the inter-rater reliability of the PRPP but this is the first investigating properties of the Norwegian version. The results of this study indicate that the inter-rater reliability of this group of occupational therapists is low compared to previous studies. The difference is, however, that this study investigated inter-reliability during the assessment course, whereas the other studies had participants that were already trained and had been using the PRPP for some time. In addition, the participants in this study used a translated version. This study supports the need for an extensive training course in order to learn the PRPP and it would not be right to expect that the occupational therapists already knew all the concepts of the PRPP on day four of the course.

This study gives some indication on the inter-rater reliability of the Norwegian version, however future studies with OT's that are well experienced with the PRPP are necessary in order to be able to draw further conclusions about the psychometric properties.

### Ethics

The participants in the study were participating on a voluntarily basis. They all signed written consent forms where it also was stated that they at any time could withdraw from the study. The clients in the videos had all given consent for the clip to be used in education and for research purposes. The consent forms were obtained by the lecturers of the PRPP training course.

### Acknowledgements

The authors would like to thank all the occupational therapists that were willing to take part in this study. It would never have been possible without all of you!

The authors would also like to thank Dr. Anders Kottorp for his valuable contributions, advice and guidance through the process working with this project.

### References

- Aubin, G; Chapparo, C, Gelinas, I; Stip, E & Rainville, C (2008). Use of the Perceive, Recall, Plan and Perform System of Task Analysis for persons with schizophrenia: A preliminary study. *Australian Occupational Therapy Journal*, Vol 56; (3) p.189-199
- American Educational Research Association (APA) (1999). Standards for Educational and Psychological Testing. American Educational Research Association, Washington
- Beaton, D.E; Bombardier, C; Gullillemin, F & Feraz, M.B (2000).

Guidelines for the process of cross-cultural adaption of self-report measures. *Spine*, 2000; 25:3186-91

Benson, J., & Clark, F. (1982) A guide for instrument development and validation. *American Journal of Occupational Therapy*. 36(12) 789-800

Birkimer, J. C. & J. H. Brown (1979). «Back to basics: Percentage agreement measures are adequate, but there are easier ways.» *Journal of applied behaviour analysis*(12): 535-543.

Chapparo, C & Ranka, J (2006). The PRPP system of Task Analysis: Users training manual

Cicerone KD, Dahlberg C, Kalmar K, Langenbahn DM, Malec JF, Bergquist TF et al. (2000). Evidence-based cognitive rehabilitation: recommendations for clinical practice. *Arch Phys Med Rehab* 2000; 81; 1596-1615.

DePoy, E & Gitlin, LN (2005). Introduction to Research: Understanding and applying multiple strategies. 3rd ed. Mosby, USA

Fisher, A (2003). Assessment of Motor and Process Skills. Development, Standardization and Administration Manual. 6th ed. Three star press, Inc

Fry, K & O`Brian, L (2002). Using the Perceive, Recall, Plan and Perform System to assess cognitive deficits in adults with traumatic brain injury: A case study. *Australian Occupational Therapy Journal*, 49, 182-187

Grieve, J & Gnanasekaran (2008). *Neuropsychology for occupational therapists. Cognition in Occupational Performance*. 3rd ed. Blackwell Publishing, UK

Laver Fawcett, A. (2008). Principles of Assessment and Outcome Measurement for occupational Therapists and Physiotherapists. Wiley & Sons, England

Law, M.C, Baum, C & Dunn, W (2005). Measuring Occupational Performance: Supporting Best Practice in Occupational Therapy. 2nd ed. SLACK Incorporated, USA

McGraw, K.O & Wong, S.P (1996). Forming inferences about some intraclass correlation coefficients. *Psychological Methods*, 1996, Vol.1, No.1, p.30-46.

Nott, M & Chapparo, C (2007). Measuring information processing in a client with extreme agitation following traumatic brain injury using the Perceive, Recall, Plan and Perform System of Task Analysis. *Australian Occupational Therapy Journal* (54)

Nott, M; Chapparo, C & Heard, R (2008). Reliability of the Perceive, Recall, Plan and Perform System of Task Analysis: A criterion-referenced assessment. *Australian Occupational Therapy Journal*.

Ottenbacher K.J. & Tomchek S.D. (1993) Reliability analysis in the-rapeutic research: practice and procedures. *American Journal of Occupational Therapy*,47:10-16.

Rankin, G. & Stokes, M. (1998). Reliability of assessment tools in rehabilitation: an illustration of appropriate statistical analysis. *Clinical rehabilitation* 1998; 12: 187-199

Shrock, S.A. & Coscarelli, W.C. (2007). Criterion-referenced test development. Technical and legal guidelines for corporate training. 3rd ed. Pfeiffer, San Francisco, USA.

Slagle, J., Weinger, M., Dinh, M.-T., Brumer, V. & Williams, K. (2002). Assessment of the intrarater and interrater reliability of an established clinical task analysis methodology. *Anesthesiology*, 96, 1129-1139.

Streiner, D. L. & G. R. Norman (2008). Health measurement scales. A practical guide to their development and use. Oxford Univeristy Press, New York.

Taylor, M.C (2007). Evidence-based Practice for Occupational Therapists. Wiley Blackwell, UK.

Thompson, B. (2003). Score reliability. Contemporary thinking on reliability issues. SAGE Publications, California, USA.

Turner, A., Foster, M. & Johnson S.E (2002). Occupational therapy and physical dysfunction. Principles, skills and practice. 5th ed. Churchill Livingstone, London, UK.

## Forsker i farten: Kjersti Vik

**Kjersti Vik er førsteamanuensis ved Høgskolen i Sør-Trøndelag og koordinator for forskningsområdet Aktivitet og Deltakelse, som er et satsningsområde ved Høgskolen i Sør-Trøndelag.**

AV SISSEL HORGHAGEN

– *Hva forsker du på?*  
 – Aktivitet og deltakelse, med hovedvekt på eldre som mottar kommunale tjenester.  
 – *Hva er dine hovedfunn?*  
 – Jeg er involvert i flere empiriske studier. Disse dokumenterer at ergoterapi og teori om aktivitet og deltakelse er viktig og nødvendig for at eldre skal ha et godt hverdagsliv. Videre viser disse studiene at eldre har ønsker om mer aktivitet og deltakelse enn de blir gitt mulighet til. Resultatet av studiene viser at det finnes mange omgivelsesfaktorer både på individ- og systemnivå som hemmer deltakelse. Videre sier de

eldre i mine studier at det viktigste som fremmer mulighet for aktivitet og deltakelse, er rullatoren og familiemedlemmene.  
 – *Hvilken forskningsartikkel er den beste du har lest, og hvorfor?*  
 – Det er det nesten umulig å svare på, men jeg vil henvise til en interessant artikkel i den Canadiske Fysioterapeut Journalen (2007, vol 59 nr 1). Her beskriver Jennifer Penney, Marilyn MacKay-Lyons, og Alison McDonald på en eminent måte hvordan man kan bruke kunnskapsbasert praksis og ICF til å dokumentere egen fagutøvelse. Artikkelen heter Evidence-Based

Stroke Rehabilitation: Case Analysis Using the International Classification of Functioning, Disability and Health Framework.  
 – *Hva trenger vi ergoterapeuter mer forskning på?*  
 – Etter hvert begynner vi å få empiri som dokumenterer ergoterapiteori. Nå må vi snart ta den vanskelige utfordringen det er å foreta studier som kan dokumentere effekt av ergoterapi. Her har vi rent metodisk noen utfordringer, og jeg kan nevne noen: Kan vi isolere faktorene som handler om effekten av ergoterapi? Hva er effekt i ergoterapi? Hvordan kan det måles?

