

Autism at a glance: A pilot study optimizing thin-slice observations

Autism

1–9

© The Author(s) 2018

Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/1362361318792872

journals.sagepub.com/home/aut

Lauren H Hampton, Philip R Curtis and Megan Y Roberts

Abstract

Borrowing from a clinical psychology observational methodology, thin-slice observations were used to assess autism characteristics in toddlers. Thin-slices are short observations taken from a longer behavior stream which are assigned ratings by multiple raters using a 5-point scale. The raters' observations are averaged together to assign a "thin-slice" value for each observation. In this study, a total of 60 toddlers were selected from a video archive: 20 children with typical development, 20 children with developmental language disorder, and 20 children with autism. In the first part of this study, 20 raters observed small play segments between toddlers and an assessor. Raters assigned scores to each of the 60 toddlers on items related to autism symptomatology. Item analysis and generalizability and decision studies were conducted to determine the factor structure and optimal number of raters to achieve a stable estimate of autism characteristics. In the second part of the study, generalizability and decision studies were conducted to determine the most efficient and optimal combination of raters and naturalistic contexts. This pilot study provides recommendations for optimizing the utility of thin-slice observations for measuring autism symptomatology in young children.

Keywords

autism spectrum disorders, screening, thin-slice

Despite the US Preventive Services Task Force conclusion that there is no clear benefit from a universal autism screening, multiple agencies, parent groups, and the American Academy of Pediatrics (AAP) all concluded that identifying autism and providing early intervention services is a public health priority and identifying efficient and effective methods for optimizing screenings is necessary (Pierce et al., 2016). Nevertheless, pediatricians do not screen for autism red flags on a routine basis (Arunyanart et al., 2012; Roux et al., 2012) or readily follow the AAP implementation recommendations, despite the fact that children with autism can be reliably identified before the age of 2 years (Cox et al., 1999; Kleinman et al., 2008; Lord et al., 2006). Families frequently wait between 1 and 3 years to receive an autism diagnosis (ADEP, 2017; Crane et al., 2016), which results in most children not receiving a diagnosis until after age 4 (Brett et al., 2016; Soke et al., 2017). Families often identify concerns about their child's communication development much younger than they receive a diagnosis (Bolton et al., 2012; Kozlowski et al., 2011). This is particularly problematic given that earlier intervention is effective for children with autism (Zwaigenbaum et al., 2015), and receiving a diagnosis is critical to gaining access to services. Children who receive an autism diagnosis are more likely to access and

benefit from early interventions and experience more positive outcomes over time (Suma et al., 2016).

Access to a skilled diagnostician is one of the greatest barriers to receiving a diagnosis (Bisgaier et al., 2011). Current screening practices rely primarily on parent report or short assessments with a novel person (Matson and Tureck, 2012; Robins et al., 2001; Stone and Ousley, 1997). The Modified Checklist for Autism in Toddlers-Revised (MCHAT-R), a widely used screening tool adapted for many different cultures, relies on parent report (Robins et al., 2001). Although this measure may be a reasonable level 1 screener, with sensitivity and specificity up to 0.90, it fails to be useful as a level 2 screener, discriminating between children with autism and children with other disabilities, such as language delays. Similarly, the Baby and Infant Screen for Children with aUtism Traits (BISCUIT; Matson and Tureck, 2012) and Screening Tool for Autism in Toddlers (STAT; Stone et al., 2004) have high sensitivity

Northwestern University, USA

Corresponding author:

Lauren H Hampton, Northwestern University, 2240 Campus Drive, Frances Searle, Room 3-365, Evanston, IL 60208, USA.

Email: lauren.hampton@northwestern.edu

and moderate specificity but again may not differentiate between autism and other disabilities well when used as the only measure. Since access to a skilled diagnostician for autism is limited (Bisgaier et al., 2011), it is critical during the time between initial screening and diagnosis (level 2) to differentiate between children who show signs of autism and need an autism-specific diagnostic testing battery, from those children who show developmental delays but no signs of autism and require a comprehensive language and developmental assessment battery. Adding a simple naturalistic direct observation to the level 2 screener stage may improve positive predictive value (PPV) and increase specificity for referrals.

Thin-slice behavior observations are short samples from a behavior stream, such as a child–assessor interaction, that contain dynamic information (Ambady and Rosenthal, 1992). Ratings are made on a series of questions by multiple raters using a 5-point scale. The ratings are averaged across raters into one thin-slice score. Within this methodology, there are three parameters that can vary: (1) the *length of segments* that raters code; (2) the *number of raters* that are averaged together to create a thin-slice score; and (3) the *number of segments watched by each rater for each child*. Although segments that are less than 5 min are recommended, the optimal length and number of observations for ratings are unknown (Carney et al., 2007). It has been suggested that as few as 5 raters or as many as 150 raters are necessary to get a useful thin-slice rating. In addition, raters from a more similar social group tend to rate more similarly than those from different social groups (Wang et al., 2009; Wright and Drinkwater, 1997; Young and Hugenberg, 2010). The simplicity of thin-slice observations may make direct observation as a screening measure more feasible if the length and number of raters are optimized.

Ratings from thin-slice observations in psychology have been used to reliably rate a variety of behaviors and are valid for predicting meaningful outcomes. These ratings have been found to be reliable on a variety of behavioral ratings of depression, Parkinson's disease, marital status, or life satisfaction (Ambady and Gray, 2002; Bonanno and Keltner, 1997; Hatfield et al., 1994; Pentland et al., 1987). Thin-slice ratings have also been used to predict important outcomes such as divorce, suicide attempts, and disease progression (Archinard et al., 2000; Gottman and Levenson, 1992; Heller and Haynal, 1997; Hertenstein et al., 2009; Kleiman and Rule, 2013; Mason et al., 2010). More recently, studies have used thin-slice observations of typically developing children to measure personality traits (Tackett et al., 2017) and to measure outcomes for children with autism (Walton and Ingersoll, 2016).

There is a clinical bottleneck in autism diagnostics, and although there are useful models across the country, the diagnostic process is still taking longer than it should and families are waiting too long for services that depend on a

diagnosis (Gordon-Lipkin et al., 2016). By optimizing the efficiency and effectiveness of the screening process, we can reduce clinical wait time and better prioritize referrals for full diagnostic evaluations. Some recent advancements in technology have led to multiple recommendations to use technology to address these needs and maximize autism discoveries. For example, crowdsourcing is one method that may accelerate our discoveries in autism (David et al., 2016). In addition, researchers have utilized content analysis of public videos to reliably detect autism (Fusaro et al., 2014), and easy-to-use mobile apps have been developed to collect videos from families for diagnostic purposes (Illingworth et al., 2017; Knutsen et al., 2016; Nazneen et al., 2015; Smith et al., 2016). Despite the potential of these innovative video-based behavioral observation methods to improve the autism diagnostic process, there is not clear evidence indicating how to best gain a stable estimate of autism symptoms from these videos. We propose that by applying the thin-slice observation methodology to these existing observational methodologies, the autism diagnostic process may be improved.

Current study

The aim of this study is to provide a preliminary assessment of the stability of multiple thin-slice ratings, provide an initial optimization of this tool for developmental screenings, and provide a set of recommendations to the field when considering thin-slice video observations as a diagnostic tool. To accomplish this aim, we conducted a two-part study to answer the following research questions: Study 1: (1) What are the dimensionality, validity, and reliability of the proposed thin-slice rating measure? Study 2: (2) What is the most efficient combination of raters and contexts to achieve a stable thin-slice value for autism symptomatology in toddlers? (3) Using an optimized thin-slice score, how does the thin-slice observation relate to similar and dissimilar measures? And finally, an exploratory question: (4) Does the optimized thin-slice value result in a meaningful PPV for identifying children with autism?

Methods

Participants

A subset of participants from a larger randomized controlled trial (RCT) was used in this retrospective study (NCT02632773). Toddlers in this larger trial were included in part by not exceeding the clinical cutoff score on the STAT at intake into the study (Stone and Ousley, 1997), yet throughout the 15-month period of the RCT study, 20 participants were identified as having autism spectrum disorder (ASD) through clinical judgment, parent report, and administration of the Autism Diagnostic Observation Scale

Table 1. Participant description.

Mean (SD)	Typical	DLD	ASD
Age in months	30.0 (4.6)	31.6	32.9
% male	90%	90%	90%
% minority	35%	20%	25%
Autism severity score	1.15 (0.49)	1.25 (0.44)	6.45 (1.88)
PLS, total language	115.95 (17.0)	74.8 (15.1)	72.6 (16.0)
Family income	\$89,236 (64,306)	\$76,336 (68,957)	\$63,501 (23,344)
NDW	86.5 (36.9)	17.0 (21.6)	37.9 (40.4)
IQ	101.25 (9.1)	93.25 (9.9)	86.5 (8.6)

SD: standard deviation; DLD: developmental language disorder; ASD: autism spectrum disorder; PLS: Preschool Language Scale, 4th edition (Zimmerman et al., 2002); NDW: number of different words from a language sample; IQ: Intelligence Quotient from the Bayley Scales of Infant and Toddler Development (Bayley, 2006).

Minority status indicates families who reported race as Black/African American, Asian, or Other and families who indicated Hispanic ethnicity. Autism severity score, higher score indicating more symptoms of autism, calculated from the Autism Diagnostic Observation Scale (Lord et al., 2008).

Table 2. Rater description.

	Study 1 (20 raters)	Study 2 (5 raters)
Mean age (SD)	24.4 (4.1)	26.4 (4.3)
Male	15%	0%
Minority	10%	0%
Some graduate school or higher	45%	60%
3 or more years' experience with toddlers	40%	80%
Experience rating videos	75%	80%

SD: standard deviation.

(ADOS; Lord et al., 2008). As such, these participants' autism symptoms were more subtle (mean ADOS severity score of 6.45, standard deviation of 1.88); this means that these children, who did not flag on a level 2 screening measure but *did* go on to receive a diagnosis of ASD, are likely the most difficult children to screen. All 20 participants diagnosed with ASD during the course of the study were included and were matched with 20 toddlers with developmental language disorder (DLD) and 20 typically developing toddlers on age, gender, and receptive language ability from the same trial. Children were matched from different single time points based on age in months, receptive/expressive language skills (for children with autism or DLD alone), and on gender. Participants are described in detail in Table 1. All assessments came from one time point within each participant, but different time points were selected across participants to improve age matching.

Raters

Raters were undergraduate or graduate speech-language pathology students or master's level research staff. Raters reported how much experience they had caring for,

observing, or working with toddler-aged children, and how much experience they had rating behavioral observations. Raters had moderate experience coding videos for language and behavior, but all were new to this method of rating and naïve to the study's purpose and to each child's diagnosis (Table 2). In Study 1, 20 raters independently rated each video using the thin-slice rating scale described in Table 3. In Study 2, a subset of five raters independently rated each video using the thin-slice rating scale across two novel observational contexts. Raters received no formal training in the measure other than definitions and examples provided in order to benefit from the thin-slice logic relying on initial impressions.

Measures

Assessments. Participants were assessed by a master's level clinician who was unfamiliar with the study purpose and hypothesis. The ADOS (Lord et al., 2008) and Preschool Language Scale, 4th edition (PLS-4; Zimmerman et al., 2002) were administered concurrently with thin-slice observations. The autism severity score, on a scale from 1 to 10, was calculated from the ADOS so that scores could be compared across modules. Scores from the ADOS are scored in a reverse manner, such that a lower score indicates a more typical performance. In addition, scores from the STAT (Stone and Ousley, 1997) were collected at entry into the study when participants were 2 years old.

A 20-min semi-structured and standardized naturalistic language sample was used to observe spontaneous communication with a novel assessor in a clinical setting. Assessors, unfamiliar with the child, played with the child using a standard set of toys (a farm set, car ramp and cars, baby doll set, play dough set, play kitchen, a book, and bubbles) and using no specific vocabulary (Miller, 1981). The adult imitated any language the child used, responded using utterances without specific language (uh-huh, wow, uhoh), and engaged in play with the child. Alternatively,

Table 3. Thin-slice average ratings.

Item	Factor loading	Typical	DLD	ASD
1. Shows interest in the adult	0.93	2.2 (0.8)	2.5 (0.5)	3.4 (0.8)
2. Plays appropriately with toys	0.80	1.9 (0.7)	2.1 (0.4)	2.5 (0.5)
3. Uses language appropriately	0.80	2.2 (0.9)	3.2 (0.7)	3.1 (0.8)
4. Uses appropriate gestures in play/communication	0.89	2.4 (0.8)	2.7 (0.7)	3.2 (0.6)
5. Imitates actions modeled by the adult	0.63	2.7 (0.8)	2.8 (0.4)	3.2 (0.3)
6. Behaves appropriately (does not engage in problem/disruptive behaviors)	0.66	1.5 (0.2)	1.8 (0.5)	1.8 (0.5)
7. Uses many different speech sounds	0.70	2.5 (1.2)	3.9 (0.9)	3.4 (1.1)
8. Understands what the adult says	0.91	2.0 (0.7)	2.5 (0.5)	2.8 (0.6)
9. Coordinates two forms of communication (eye contact, gestures, vocalizations, or words)	0.90	2.2 (0.9)	2.6 (0.8)	3.2 (0.9)
10. Uses appropriate eye contact	0.82	2.3 (0.8)	2.5 (0.7)	3.4 (0.7)
11. Engages in unusual behaviors	0.90	2.1 (0.5)	2.5 (0.6)	3.0 (0.4)
Total thin-slice score		24.2 (7.7)	28.6 (5.5)	33.0 (6.6)

DLD: developmental language delay; ASD: autism spectrum disorder.

All items are scored such that lower scores are closer to typical development and higher scores indicate greater impairment.

parent–child interactions included a different standard set of toys (doll house, ball, school bus set, book, puzzle, and blocks) and parents were instructed to play and communicate with their child as they usually do (Roberts and Kaiser, 2012). Language samples and parent–child interactions were video recorded and transcribed and coded for the number of different words (NDW) spontaneously produced by the child. Fidelity of administration and verification of standardized scoring was assessed on 20% of administrations and exceeded 90% across all measures. And 20% of the language samples were coded by a second coder, and point-by-point agreement exceeded 90% across language samples.

Thin slices. Thin-slice observations were created by selecting a 2-min segment from each of the three contexts: the play portion of the ADOS modules 1 and 2, a naturalistic language sample conducted by an unfamiliar assessor that mirrors the play portion of the ADOS, and a parent–child, play-based interaction. Three different contexts were selected to retain independence of the samples; the first interaction was selected to validate the rating scale and the second two interactions were selected to determine stability and validity of thin-slice scores across the contexts and for the greatest applicability as a screening measure. Segments of 2-mins were selected starting at 3 min into the original video recording, or the next 2 min of uninterrupted free-play, whichever came first. Interruptions included 5 s or longer of the child or adult being off camera or segments where the adult was speaking to someone other than the child.

Raters were randomly assigned a unique assignment order such that no two raters observed the videos in the same sequence. The raters viewed each 2-min segment without pausing or rewinding the video. Raters assigned codes for each video on 11 unique items, rated on a 5-point

Likert-type scale (see Table 3). These ratings from all 11 items were averaged together so that each rater had one thin-slice score for each video. In Study 1, raters observed the thin slice from the ADOS; in Study 2, the raters observed the language sample and parent interaction videos intermixed together. Each item in the 11-item thin-slice rating scale was defined and raters were instructed to rate the quality of the skill or behavior they observed in the child on the anchored scale from 1 to 5 (Table 3). The raters were also instructed to use clinical judgment based on experiences with typically developing 2- to 3-year-olds.

Thin-slice rating items. The 11 unique items were derived from multiple sources for this preliminary study, of which we considered 10. First, the 5-item scale used by Walton and Ingersoll (2016) to rate communicative behaviors in children with ASD were replicated (items 1–5). Second, additional behaviors that may be relevant to discriminate children with DLD from children with ASD were added (items 6–8). Third, three items were created using unique variables from the ADOS and reworded to fit with the other items in the scale (items 9–11). Raters were instructed to rate each item based on how much they agreed that the statement was either characteristic of the child in the short clip (1=strongly disagree/never observed, 2=disagree/observed <5% of observation, 3=neither agree nor disagree/observed <50% of observation, 4=agree/observed >50% of observation, 5=strongly agree and >80% of observation).

Analysis

Item reliability and validity. To determine the structure of the 11 items used to code the thin-slices, average ratings from 20 raters on 60 slices from the ADOS were used in a factor analysis. All analyses were run in R using the psych

(Revelle, 2011) and paran packages (Dinno, 2009). In order to determine the number of factors within the scale, Horn's Parallel Analysis was used and results supported the retention of only one factor (Horn, 1965). Likewise, a Very Simple Structure analysis also suggested that only one factor was supported from the current data (Revelle and Rocklin, 1979). A factor analysis was run to extract one factor, estimated using the minimum residual method. The solution was subjected to an oblimin rotation. All items loaded highly on this one extracted factor (factor loadings are given in Table 3). These results support the unidimensionality of the items used in coding the thin-slices. Reliability of the scale was assessed using both Cronbach's α as well as the more conservative omega-hierarchical coefficient (ω_h ; Zinbarg et al., 2005). Omega-hierarchical uses a hierarchical factor model, in which both a general factor and lower order factors can account for variance in observed scores; the ω_h coefficient represents the general factor saturation.

G-study. Generalizability theory (G-theory) allows us to determine how much variability across facets (e.g. observations, contexts, observers) of a study is due to true score relative to error (Cronbach et al., 1972). A G-coefficient uses the partitioned variance attributed to each measurement facet, similar to an interclass correlation, to give an estimate of measurement stability. Specifically, the G-coefficient is calculated using components from a traditional analysis of variance (ANOVA; equation (1)). This ratio separates variance due to the individual person from the total combined variance components of the person, person by rater, person by session, and total combined error. Increased stability in measurement increases confidence that the measure is representative and therefore more likely generalizable, thus increasing external validity. A G-coefficient between 0.6 and 0.8 is considered an acceptable level of stability (Bakeman et al., 1997).

D-study. While a G-study results in an estimate of achieved stability, a determination study (D-study) is used to estimate the expected achieved G-coefficient with increased samples (Shavelson et al., 1989; equation (2)). A D-study uses the fixed variance from the G-study and applies it to additional samples, and therefore assumes that the variability attributed to the original samples remains stable over time. The relative error variance (σ_{rel}^2) is a division of the individual variances by different combinations of hypothesized sample sizes for a number of sessions and raters (see formula (2) and the calculator available from Webb and Shavelson (2005)). A D-study allows the researcher to optimize a measure by estimating the ideal measurement samples necessary to achieve an optimal G-coefficient. For example, when considering how many different raters are needed to achieve a stable G-coefficient, a D-study may estimate two raters at $g=0.59$, three raters at $g=0.72$,

and four raters at $g=0.80$. Since three raters and four raters both fall in the recommended range, the research can select to average three rater observations together to optimize resources and outcomes. All G- and D-studies were completed using the EduG (2012) software.

Equations

$$G = \frac{\sigma_c^2}{\sigma_c^2 + \sigma_{cr}^2 + \sigma_{cs}^2 + \sigma_{crs,e}^2} \quad (1)$$

$$Dstudy: \sigma_{rel}^2 = \left(\frac{\sigma_{cr}^2}{n_r'} \right) + \left(\frac{\sigma_{cs}^2}{n_s'} \right) + \left(\frac{\sigma_{crs,e}^2}{n_r' n_s'} \right) \quad (2)$$

where σ^2 represents the variance components, c the child-level variance components, r the raters, s the sessions, e the error variance, and n is the sample size (Yoder and Symons, 2010: 29).

Measurement validity. Once an optimal G-coefficient was identified, the optimal numbers of raters and contexts of the thin-slice were averaged together to result in an optimized thin-slice score. Correlations between thin-slice scores and measures of similar constructs (language and ASD measures) were estimated to demonstrate validity. In addition, thin-slice scores were correlated with divergent variables (the child's age and socioeconomic status (SES)) expected to not be highly correlated with the outcome to further increase validity of the measure.

Screening utility. To estimate the utility of the optimized thin-slice score as a screening tool, the PPV and negative predictive value (NPV) were calculated to estimate the utility of the thin-slice measure in correctly categorizing toddlers with and without autism as compared to children with DLD. Sensitivity, specificity, NPV, and PPV were also calculated across all three populations to compare the utility to other screeners. These likelihood ratios consider the proportion of accurately classified children with autism to those accurately classified as not having autism. A receiver operating characteristic (ROC) curve analysis, with bootstrapped confidence intervals, was used to estimate the best cut score for optimizing sensitivity, specificity, NPV, and PPV using pROC package in R (Robin et al., 2011).

Results

Study 1

To determine the structure of the 11 items used to code the thin slices, average ratings from 20 raters on 60 slices from the ADOS were used in a factor analysis. All analyses were run in R using the psych (Revelle, 2011) and paran

Table 4. Correlations between measures and thin-slice score.

	Age (months)	Family income	ADOS: severity ^a	PLS: total standard	NDW	Total slice score
Age (months)	I					
Family income	-0.03	I				
ADOS: severity ^a	0.03	-0.09	I			
PLS: total standard	-0.10	0.21	-0.45*	I		
NDW	0.35*	0.16	-0.23	0.77*	I	
Total slice score	-0.07	-0.09	0.47*	-0.59*	-0.53*	I

ADOS: Autism Diagnostic Observation Scale (Lord et al., 2008); PLS: Preschool Language Scale, 4th edition (Zimmerman et al., 2002); NDW: number of different words.

^aHigher scores indicate greater autism severity.

* $p < 0.05$.

packages (Dinno, 2009). In order to determine the number of factors within the scale, Horn's Parallel Analysis was used and the results supported the retention of only one factor (Horn, 1965). Likewise, a Very Simple Structure analysis also suggested that only one factor was supported from the current data (Revelle and Rocklin, 1979). When a factor analysis was run with one factor, all items loaded highly on this one factor (all factor loadings > 0.62). These results support the unidimensionality of the items used in coding the thin-slices.

Reliability of the scale was assessed using both Cronbach's α as well as the more conservative omega-hierarchical coefficient (ω_h ; Zinbarg et al., 2005). Omega-hierarchical models a hierarchical factor model, in which both a general and item-specific factors can account for variance in observed scores; the ω_h coefficient represents the degree to which items load onto the *general factor*. Both α and ω_h estimates suggested excellent reliability ($\alpha = 0.95$, $\omega_h = 0.82$). In order to test the validity of the new scale in assessing autism symptomatology, average sum scores across raters from the 60 thin slices were correlated with the STAT at baseline ($r = 0.57$) and time-point-specific ADOS severity scores ($r = 0.35$). Correlations were moderate and significant and support the validity of the new scale.

Using the initial 60 ratings from the ADOS, a single-facet G-study was conducted to optimize the number of raters required to establish a stable estimate of the thin-slice average score for each child. The single-facet G-study resulted in an absolute G-coefficient of 0.95 when using 20 raters. Interestingly, the D-study, which uses a bootstrap-like resampling method to generalize to larger samples, found that an optimal rating could be achieved with as few as two ($g = 0.69$) and five ($g = 0.82$) raters of a semi-structured observation. Due to the subjectivity of the rating and the potential for increased variability when more observations are introduced, five raters were selected for Study 2 to ensure maximal optimization.

Study 2

Optimized across contexts. The two-facet G-study, used to determine the optimal number of contexts and raters, resulted

Table 5. Number of participants flagged as ASD risk using 34.3 cut score.

	Slice rating		
	ASD risk	No ASD risk	
Clinical diagnosis	ASD	16	4
	DLD	7	13
	Typical	0	20

ASD: autism spectrum disorder, DLD: developmental language delay.

in an optimal stable estimate when five raters rated two novel naturalistic contexts ($g = 0.73$). This was a large improvement over the G-coefficient obtained when five raters rated only one naturalistic context ($g = 0.60$). During the second study, two new contexts were selected to reduce bias, yet greater instability was observed when five raters were used as compared to Study 1. Therefore, we selected the most optimized measure that promised the greatest generalizability across naturalistic contexts and averaged the five raters' scores across both observations (language sample and parent-child interaction), using the optimized thin-slice score for the subsequent analyses. In addition, the language sample and parent-child interaction have the most utility as a screening measure as opposed to the ADOS thin slice.

Correlations. The correlations between the optimized thin-slice score and related measures of language and autism were moderate ($r = -0.58$, -0.47 , $p < 0.05$; see Table 4). Even though NDW also significantly related to the optimized thin-slice score ($r = -0.59$, $p < 0.05$), NDW and autism severity were not significantly correlated with one another ($r = -0.23$, $p > 0.1$). Divergent scores (SES and age) were not significantly or meaningfully correlated with the optimized thin-slice score ($r = -0.11$, -0.02 , $p > 0.1$).

PPV. An ROC curve analysis was used to estimate an optimized cut score of 34.3 (AUC = 0.82; Robin et al., 2011; Table 5). This cut score resulted in high sensitivity and specificity for distinguishing autism from DLD and typical development (Table 6). In addition, moderate PPV and

Table 6. Diagnostic accuracy (95% bootstrapped confidence interval).

	Sensitivity	Specificity	PPV	NPV
ASD versus DLD	0.80 [0.60–0.95]	0.65 [0.45–0.85]	0.70 [0.57–0.84]	0.76 [0.60–0.94]
ASD versus DLD + typical	0.80 [0.60–0.95]	0.83 [0.70–0.93]	0.70 [0.56–0.86]	0.89 [0.81–0.97]

ASD: autism spectrum disorder; DLD: developmental language delay; PPV: positive predictive value; NPV: negative predictive value. 95% confidence intervals bootstrapped with 10,000 stratified replicates.

NPV resulted from identifying children with autism from those with DLD alone (see Table 6). Thus, the proposed thin-slice rating tool may be a useful tool toward estimating autism in young children; however, this may be further optimized when combined with additional level 2 screening measures.

Discussion

The thin-slice measure resulted in stable estimates across a relatively few number of raters and contexts. Using short observations of naturalistic interactions between young children and their parents and a short interaction with an assessor may provide a useful perspective on the child's development and risk for autism. Optimized thin-slice ratings were significantly correlated with the severity score from the ADOS, the gold standard (but time-intensive) assessment tool in autism diagnostics, and the STAT measure for screening. Unrelated variables were not correlated with the thin-slice score (age, family income). However, the relationship between the thin-slice scores and language measures may indicate that the thin-slice scores may over-emphasize language and communication skills. Finally, the optimized thin-slice score was relatively effective in identifying children with autism as compared to children with DLD (PPV=0.70; NPV=0.83) in this preliminary study. This PPV is interesting given that the children in this study had lower autism severity scores and were compared to children with DLD, a comparison that is often difficult for differential diagnosis (Bishop and Norbury, 2002; Loucas et al., 2008). Although this PPV is not high enough to recommend using the thin-slice score as a screener alone, it suggests that the thin-slice score could add meaningful value to a battery of screening assessments for children suspected of having autism.

Implications

This study is an important first step toward optimizing the thin-slice measure as a useful tool for characterizing autism symptomatology in young children. One of the important results of this study is the practicality of utilizing the thin-slice methodology in clinical practice for toddlers with autism. Only five raters' scores on two 2-min contexts (4 min per coder, 20 min of total coding time across all raters) were needed to arrive at a stable estimate of children's social communication scores. This means that utilizing the thin-slice

scoring methodology may be an efficient and cost-effective measurement tool. Although the thin slice is an interesting new measure, it is important to understand how to best optimize and utilize a measure that relies on short observations to maximize generalizability.

Limitations

The results of this study must be considered in light of three primary limitations. First, the ratings were conducted by well-trained research staff who are familiar with rating and coding child communicative behavior. The results of the G- and D-studies may only generalize to other raters with experience working with young children. Second, the ratings during Study 2 came from different contexts than Study 1. Although all contexts were play based, the variability in G-coefficients across contexts suggests that further replication is necessary to provide a truly optimized measure. Third, this study used a video archive and was limited by the participants available. A prospective screening study is necessary to replicate this work.

Next steps

Thin-slice rating requires the use of a multi-item, Likert-type scale. This study used a novel, 11-item scale for this purpose. The scale showed strong reliability and acceptable validity within the current sample. The 11 items in this study, though not exhaustive, were selected to replicate previous work using thin-slice measures and to provide a context for optimizing the thin-slice approach. Further expansion and refinement of the 11-item scale may result in improved predictive validity. In addition, to further validate the thin-slice observation scale, a predictive validity study should be conducted. In addition, the G- and D-studies should be replicated to increase the optimization and generalizability across populations and contexts. The accuracy of ratings for less experienced raters should also be examined in order to determine the needed qualifications and experience for reliable thin-slice coding in clinical practice. Future studies should consider the thin-slice measure along with existing screening tools and the additive value of this observational measure. Finally, the thin-slice measure should be explored for other important uses, such as an outcome measure for intervention studies aiming to reduce autism symptomatology and increase social communication.

Conclusion

This study provides preliminary support for the use of thin-slice coding procedures using short, 2-min video-taped interactions in order to assess children's social communication and autism symptomatology. Short segments from a behavioral stream, or thin-slices, can be optimized and reliability rated when averaged across five raters and two contexts. The thin-slice observation is an exciting and useful tool for making snap judgment that may be effective for rating autism characteristics. Although additional validation is needed, this short observational measure may be an effective tool to increase efficiency in autism screenings.

Funding

The author(s) received no financial support for the research, authorship, and/or publication of this article.

References

- ADEP (2017) Autism diagnosis education project. *Ohio Center for Autism and Low Incidence*. Available at: <http://www.ocali.org/project/adep> (accessed 14 November 2017).
- Ambady N and Gray HM (2002) On being sad and mistaken: mood effects on the accuracy of thin-slice judgments. *Journal of Personality and Social Psychology* 83(4):947–961.
- Ambady N and Rosenthal R (1992) Thin slices of expressive behavior as predictors of interpersonal consequences: a meta-analysis. *Psychological Bulletin* 111: 256–274.
- Archinard M, Haynal-Reymond V and Heller M (2000) Doctor's and patients' facial expressions and suicide reattempt risk assessment. *Journal of Psychiatric Research* 34(3):261–262.
- Arunyanart W, Fenick A, Ukritchon S, et al. (2012). Developmental and autism screening: A survey across six states. *Infants & Young Children* 25(3): 175–187.
- Bakeman R, McArthur D, Quera V, et al. (1997) Detecting sequential patterns and determining their reliability with fallible observers. *Psychological Methods* 2(4): 357–370.
- Bayley N (2006) *Bayley Scales of Infant and Toddler Development*. London: Pearson.
- Bisgaier J, Levinson D, Cutts DB, et al. (2011) Access to autism evaluation appointments with developmental-behavioral and neurodevelopmental subspecialists. *Archives of Pediatrics & Adolescent Medicine* 165(7): 673–674.
- Bishop DVM and Norbury CF (2002) Exploring the borderlands of autistic disorder and specific language impairment: a study using standardised diagnostic instruments. *Journal of Child Psychology and Psychiatry* 43(7): 917–929.
- Bolton PF, Golding J, Emond A, et al. (2012) Autism spectrum disorder and autistic traits in the Avon Longitudinal Study of Parents and Children: precursors and early signs. *Journal of the American Academy of Child and Adolescent Psychiatry* 51(3): 249–260.
- Bonanno GA and Keltner D (1997) Facial expressions of emotion and the course of conjugal bereavement. *Journal of Abnormal Psychology* 106(1): 126–137.
- Brett D, Warnell F, McConachie H, et al. (2016) Factors affecting age at ASD diagnosis in UK: no evidence that diagnosis age has decreased between 2004 and 2014. *Journal of Autism and Developmental Disorders* 46(6): 1974–1984.
- Carney DR, Colvin CR and Hall JA (2007) A thin slice perspective on the accuracy of first impressions. *Journal of Research in Personality* 41(5): 1054–1072.
- Cox A, Klein K, Charman T, et al. (1999) Autism spectrum disorders at 20 and 42 months of age: stability of clinical and ADI-R diagnosis. *Journal of Child Psychology and Psychiatry and Allied Disciplines* 40(5): 719–732.
- Crane L, Chester JW, Goddard L, et al. (2016) Experiences of autism diagnosis: a survey of over 1000 parents in the United Kingdom. *Autism* 20(2): 153–162.
- Cronbach LJ, Gleser GC, Nanda H, et al. (1972) *The Dependability of Behavioral Measurements: Theory of Generalizability for Scores and Profiles*. New York: John Wiley & Sons.
- David MM, Babineau BA and Wall DP (2016) Can we accelerate autism discoveries through crowdsourcing? *Research in Autism Spectrum Disorders* 32(Suppl. C): 80–83.
- Dimino A (2012) paran: Horn's test of principal components/factors. R package version 1.5.1.
- EduG (2012) EduG (version 6.1-e, generalizability study) [Société Suisse pour la Recherche en Éducation, Groupe de travail Edumétrie—Qualité de l'évaluation en éducation]. Longueuil, QC, Canada: Software Prepared by Maurice Dalois and Léo Laroche.
- Fusaro M, Vallotton CD and Harris PL (2014) Beside the point: mothers' head nodding and shaking gestures during parent-child play. *Infant Behavior & Development* 37(2): 235–247.
- Gordon-Lipkin E, Foster J and Peacock G (2016) Whittling down the wait time: exploring models to minimize the delay from initial concern to diagnosis and treatment of autism spectrum disorder. *Pediatric Clinics of North America* 63(5): 851–859.
- Gottman JM and Levenson RW (1992) Marital processes predictive of later dissolution: behavior, physiology, and health. *Journal of Personality and Social Psychology* 63(2): 221–233.
- Hatfield E, Cacioppo JT and Rapson RL (1994) *Emotional Contagion: Cambridge Studies in Emotion and Social Interaction*. Cambridge: Cambridge University Press.
- Heller M and Haynal V (1997) Depression and suicide faces. In: Ekman P and Rosenberg EL (eds) *What the Face Reveals*. Oxford: Oxford University Press, pp.398–407.
- Hertenstein MJ, Hansel CA, Butts AM, et al. (2009) Smile intensity in photographs predicts divorce later in life. *Motivation and Emotion* 33(2): 99–105.
- Horn JL (1965) A rationale and test for the number of factors in factor analysis. *Psychometrika* 30(2): 179–185.
- Illingworth DA, Thomas RP, Rozga A, et al. (2017) Cue use in distal autism spectrum assessment: a lens model analysis of the efficacy of telehealth technologies. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 61(1): 170.
- Kleiman S and Rule NO (2013) Detecting suicidality from facial appearance. *Social Psychological and Personality Science* 4(4): 453–460.
- Kleinman JM, Ventola PE, Pandey J, et al. (2008) Diagnostic stability in very young children with autism spectrum disorders. *Journal of Autism and Developmental Disorders* 38(4): 606–615.

- Knutsen J, Wolfe A, Burke BL, et al. (2016) A systematic review of telemedicine in autism spectrum disorders. *Review Journal of Autism and Developmental Disorders* 3(4): 330–344.
- Kozlowski AM, Matson JL, Horovitz M, et al. (2011) Parents' first concerns of their child's development in toddlers with autism spectrum disorders. *Developmental Neurorehabilitation* 14(2): 72–78.
- Lord C, Risi S, DiLavore PS, et al. (2006) Autism from 2 to 9 years of age. *Archives of General Psychiatry* 63(6): 694–701.
- Lord C, Rutter M, Dilavore PC, et al. (2008) *ADOS: Autism Diagnostic Observation Schedule*. Boston, MA: Hogrefe.
- Loucas T, Charman T, Pickles A, et al. (2008) Autistic symptomatology and language ability in autism spectrum disorder and specific language impairment. *Journal of Child Psychology and Psychiatry* 49(11): 1184–1192.
- Mason AE, Sbarra DA and Mehl MR (2010) Thin-slicing divorce: thirty seconds of information predict changes in psychological adjustment over 90 days. *Psychological Science* 21(10): 1420–1422.
- Matson JL and Tureck K (2012) Early diagnosis of autism: current status of the Baby and Infant Screen for Children with aUtism Traits (BISCUIT-Parts 1, 2, and 3). *Research in Autism Spectrum Disorders* 6(3): 1135–1141.
- Miller JF (1981) *Assessing Language Production in Children: Experimental Procedures*, vol. 1. Baltimore, MD: University Park Press.
- Nazneen N, Rozga A, Smith CJ, et al. (2015) A novel system for supporting autism diagnosis using home videos: iterative development and evaluation of system design. *JMIR mHealth and uHealth* 3(2): e68.
- Pentland B, Pitcairn TK, Gray JM, et al. (1987) The effects of reduced expression in Parkinson's disease on impression formation by health professionals. *Clinical Rehabilitation* 1(4): 307–312.
- Pierce K, Courchesne E and Bacon E (2016) To screen or not to screen universally for autism is not the question: why the task force got it wrong. *Journal of Pediatrics* 176: 182–194.
- Revelle W (2011) Psych: Procedures for personality and psychological research. Evanston, IL: Northwestern University. Available at: <http://personality-project.org/r/psych.manual.pdf> (accessed 1 January 2009).
- Revelle W and Rocklin T (1979) Very Simple Structure: an alternative procedure for estimating the optimal number of interpretable factors. *Multivariate Behavioral Research* 14: 403–414.
- Roberts MY and Kaiser AP (2012) Assessing the effects of a parent-implemented language intervention for children with language impairments using empirical benchmarks: a pilot study. *Journal of Speech, Language, and Hearing Research* 55(6): 1655–1670.
- Robin X, Turck N, Hainard A, et al. (2011) pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinformatics* 12(1): 77.
- Robins DL, Fein D, Barton ML, et al. (2001) The Modified Checklist for Autism in Toddlers: an initial study investigating the early detection of autism and pervasive developmental disorders. *Journal of Autism and Developmental Disorders* 31(2): 131–144.
- Roux AM, Herrera P, Wold CM, et al. (2012) Developmental and autism screening through 2-1-1. *American Journal of Preventive Medicine* 43(Suppl. 5): S457–S463.
- Shavelson RJ, Webb NM and Rowley GL (1989) Generalizability theory. *American Psychologist* 44(6): 922–932.
- Smith CJ, Rozga A, Matthews N, et al. (2016) Investigating the accuracy of a novel telehealth diagnostic approach for autism spectrum disorder. *Psychological Assessment* 29: 245–252.
- Soke GN, Maenner MJ, Christensen D, et al. (2017) Brief report: estimated prevalence of a community diagnosis of autism spectrum disorder by age 4 years in children from selected areas in the United States in 2010: evaluation of birth cohort effects. *Journal of Autism and Developmental Disorders* 47(6): 1917–1922.
- Stone WL and Ousley OY (1997) STAT manual: screening tool for autism in two-year-olds. Unpublished manuscript, Vanderbilt University, Nashville, TN.
- Stone WL, Coonrod EE, Turner LM, et al. (2004) Psychometric properties of the STAT for early autism screening. *Journal of Autism and Developmental Disorders* 34(6): 691–701.
- Suma K, Adamson LB, Bakeman R, et al. (2016) After early autism diagnosis: changes in intervention and parent-child interaction. *Journal of Autism and Developmental Disorders* 46(8): 2720–2733.
- Tackett JL, Smack AJ, Herzhoff K, et al. (2017) Measuring child personality when child personality was not measured: application of a thin-slice approach. *Personality and Mental Health* 11(1): 4–13.
- Walton KM and Ingersoll BR (2016) The utility of Thin Slice ratings for predicting language growth in children with autism spectrum disorder. *Autism* 20(3): 374–380.
- Wang EJ, Toosi NR and Ambady N (2009) Nonverbal dialects: culture and person perception. In: Wyer RS, Chiu CY, Hong YY, et al. (eds) *Understanding Culture: Theory, Research and Application*. New York: Psychology Press, pp.289–298.
- Webb NM and Shavelson RJ (2005) Generalizability theory: overview. In: Everitt BS and Howell DC (eds) *Encyclopedia of Statistics in Behavioral Science*, vol. 2. Chichester: John Wiley & Sons, pp.717–719.
- Wright JC and Drinkwater M (1997) Rationality vs accuracy of social judgment. *Social Cognition* 15(4): 245–273.
- Yoder P and Symons F (2010) *Observational Measurement of Behavior*. New York: Springer Publishing Company.
- Young SG and Hugenberg K (2010) Mere social categorization modulates identification of facial expressions of emotion. *Journal of Personality and Social Psychology* 99(6): 964–977.
- Zimmerman IL, Steiner VG and Pond RE (2002) *PLS-4: Preschool Language Scale*. 4th ed. San Antonio, TX: The Psychological Corporation.
- Zinbarg RE, Revelle W, Yovel I, et al. (2005) Cronbach's α , Revelle's β , and McDonald's ω H: their relations with each other and two alternative conceptualizations of reliability. *Psychometrika* 70: 123–133.
- Zwaigenbaum L, Bauman ML, Choueiri R, et al. (2015) Early intervention for children with autism spectrum disorder under 3 years of age: recommendations for practice and research. *Pediatrics* 136(Suppl. 1): S60–S81.