TECHNICAL APPENDIX TO THE ARRT'S ANNUAL REPORT OF EXAMINATIONS: PRIMARY ELIGIBILITY PATHWAY RESULTS 2019



Introduction

This report summarizes the psychometric characteristics of ARRT's examination scores in Radiography (RAD), Nuclear Medicine Technology (NMT), Radiation Therapy (THR), Sonography (SON), and Magnetic Resonance Imaging (MRI) for the year 2019. This report is a companion document to the *Annual Report of Examinations: Primary Eligibility Pathway* report.

The first section of this report contains information about the duration of time that candidates used to complete their examinations. The second section provides descriptive statistics of total exam scores, both raw and scaled, and information about how ARRT converts raw scores to scaled scores. The third section of this report presents descriptive statistics for the exams' section scores, including correlations and reliability estimates. Section four provides more detail about the reliability of the overall exam scores, with a discussion of coefficient α and the standard error of measurement. The final section of the report addresses decision consistency, which quantifies the reproducibility of the certification and registration decisions that ARRT makes based on its examinations.

Information about Exam Durations

Most examination administrators, including ARRT, do not intend to have exam administration time be a heavily influential factor for candidates. Practical limitations, however, make it necessary to establish exam time limits. For RAD, NMT, THR, and MRI, candidates may take up to 210 minutes (3.5 hours) to answer 220 items (questions). For SON, candidates may take up to 390 minutes (6.5 hours) to answer 400 items. The intention of the time limit is to have the exam begin and end in a reasonable amount of time, while also ensuring that knowledgeable candidates have sufficient time to complete the exam assuming that they remain focused. It is ARRT's intention that, although its exams are time limited, its exams are not speeded exams.

This section presents information on the amount of time that candidates used to take the exams described in this report. Some sources (e.g., Nunnally, 1978) specify that an exam is unspeeded when at least 90% of candidates complete the exam within the allotted time. If results show that more than 10% of candidates require the full time, ARRT would consider re-evaluating existing time limits.

Table 1 contains a summary of the amount of exam time candidates spent. These and all other statistics reflect only first-time ARRT exam candidates. None of the statistics include state candidates or people retaking the exam after failing the initial attempt. This table indicates that THR candidates spent more time than their counterparts in RAD, NMT, and MRI. THR had the highest mean (average) time among the exams with 200 questions. SON took more time overall, but the time per item was lower than the other four disciplines.

Table 1. Dest	I puve Diatistic	.s of Canaluates	The Spent of		(m minuces)
Discipline	Number of	Minimum	Maximum	Mean	Standard
Discipline	Candidates	Time	Time	Time	Deviation
RAD	11,769	38	210	142.20	38.67
NMT	394	49	210	137.70	39.84
THR	823	69	210	160.30	35.22
SON	397	99	389	229.47	70.97
MRI	515	54	210	145.50	41.95
	515	51	210	115.50	11.95

Table 1. Descriptive Statistics of Candidates' Time Spent on Examination (in Minutes)

Table 2 divides the candidates into nine groups according to the amount of time for the cumulative group to complete the exam. Using RAD as an example, 10% of all candidates completed the exam in 92 minutes or less, and 20% completed it in 106 minutes or less. Continuing on the row, Table 2 shows that 90% of RAD candidates completed the exam in 199 minutes or less. Overall, most candidates completed their examinations within the established time limits. For all disciplines, 90% or more of the candidates completed the exam in less than the allotted time. These exams do not appear to be speeded under the 90% or more completion criterion.

140	Table 2. Number of Minutes Required to Complete Exams by Tereentiles								
Dissipling		Cumu	lative Per	centage o	f Candida	tes Comp	leting the	Exam	
Discipline -	10%	20%	30%	40%	50%	60%	70%	80%	90%
RAD	92	106	117	129	140	152	166	181	199
NMT	85	99	112	123	136	149	159	179	196
THR	110	127	139	151	162	174	184	197	207
SON	142	161	183	203	220	240	259	293	340
MRI	87	103	119	134	148	162	175	190	201

Table 2. Number of Minutes Required to Complete Exams by Percentiles

Descriptive Statistics for Total Examination Scores

Table 3 contains descriptive statistics for the raw scores (number correct), which are the basis for numerous other calculations in this report. The total score consists of 200 items for RAD, NMT, THR, and MRI. The total score consists of 360 items for SON. There are also additional unscored "pilot" items on each exam.

	Table 3. Descriptive Statistics of Raw Scores						
Discipline	Minimum	Maximum	Mean	Standard			
				Deviation			
RAD	50	198	157.36	19.48			
NMT	85	189	149.45	21.21			
THR	84	194	155.30	17.54			
SON	121	341	258.93	43.38			
MRI	57	192	148.17	25.72			



ARRT uses scaled scores to report exam results. Total scaled scores range from 1 to 99, and a candidate must achieve a total scaled score of 75 to pass an examination. Table 4 contains descriptive statistics for the total scaled scores. The main advantage of scaled scores is that they facilitate a meaningful comparison of scores across forms and years.

Discipline	Minimum	Maximum	Mean	Standard Deviation
RAD	43	99	83.36	7.37
NMT	62	96	83.20	6.83
THR	54	97	81.86	6.95
SON	50	96	78.54	8.99
MRI	46	97	80.33	9.52

In order to convert raw scores to scaled scores, ARRT determines the difficulty of an exam form. Each exam consists of items that were used on previous exams. ARRT uses the Rasch model to track the difficulty levels of individual exam items and, consequently, whole exam forms. Each item has a Rasch difficulty statistic indicating the probability of a candidate answering correctly.

ARRT determines the difficulty of an exam form by calculating the sum of the probabilities of correct answers at the cutpoint. Comparisons with the difficulties of previous forms determine the relative difficulty level of the new form. If the new form is easier, the cut score for the new form will be greater by an appropriate number of questions. If the new form is more difficult, then the cut score will be lower by some appropriate number of questions.

After determining the raw passing score, ARRT calculates equations to convert the raw scores to scaled scores such that the scaled scores range from 1 to 99 with a passing score of 75. As a hypothetical example, assume that the raw passing score is 130 out of 200. The conversion equation requires two scaling coefficients: the slope (a) and the intercept (b). The calculations of a and b involve four values: the maximum scaled score (99.49), the scaled cut score (74.50), the maximum raw score (200), and the raw cut score (130).

$$a = (99.49 - 74.50) / (200 - 130) = 0.357$$

 $b = 74.50 - (a \times 130) = 74.50 - (0.357 \times 130) = 28.09$

For this hypothetical form, the scaling coefficients would be a = 0.357 and b = 28.09. ARRT would use these scaling coefficients to convert the raw scores to scaled scores. If a candidate achieved a raw score of 131 (one point above passing), then the scaled score would be

scaled score = (raw score
$$\times 0.357$$
) + 28.09 = (131 $\times 0.357$) + 28.09 = 74.857,

which rounds up to 75, a passing scaled score. For this example, raw scores of 130 and 131 round up to a passing scaled score of 75. Raw scores of 128 and 129, however, round to a scaled score of 74, which is a failing score.



Table 5 contains the pass percentages for exams taken by primary pathway candidates. This information is also in the *Annual Report of Examinations: Primary Eligibility Pathway* report but is repeated here because of its importance. One item of note is that on January 1, 2019, ARRT began using a new cut score for the Radiation Therapy examination. Although ARRT continued to report the exam's cut score as a scaled score of 75, the new cut score required passing candidates to answer a few more questions correctly than they had to in the past.

Pass Percentage
88.98
89.59
86.63
63.48
76.70

Table 5. Pass Percentages for First-Time Candidates

Descriptive Statistics for Section Scores

In addition to the total scaled score, ARRT reports individual section scores that correspond to content areas as outlined in the content specifications of each exam. The primary purpose of the section scores is to provide general information to candidates regarding their strengths and weaknesses in particular content categories. For SON only, candidates must pass both the Abdomen and OB/GYN sections in addition to passing the exam as a whole. ARRT reports section scores on a scale from 0.1 to 9.9 in one-tenth point intervals.

Section scores are useful to the extent that: (a) the scores are reliable and (b) the sections measure knowledge and skills that are independent of each other. For these reasons, Tables 6 through 10 contain additional descriptive statistics about ARRT's section scores. These include the correlations among the section scores as well as the number of items in each section, raw score means, and standard deviations. In addition, the tables contain a reliability estimate (Cronbach's α) for each section. Sections with more items generally have more reliable scores in the same way that longer examinations generally have more reliable scores. Reliability is discussed in more detail later in this report.

The correlations among the section scores provide a measure of their distinctness. In theory, correlations can range from -1.00 (perfect inverse linear relationship) to +1.00 (perfect positive linear relationship). Section scores on an exam are usually positively correlated, because candidates who perform well on one section typically perform well on others. In Tables 6 through 10, the section score correlations above the diagonal are the observed (uncorrected) correlations, and the correlations below the diagonal are correlations corrected for unreliability. The corrected correlations take into account the unreliability of the section scores and give a sense of the magnitude of the correlations under the condition of perfect reliability. The high correlations after correction among many of the section scaled scores indicate a high degree of common variance among these scores.



	I dole 0			correlati			- CICD	
Content Area	PC1	S 1	S2	IP1	IP2	P1	P2	P3
PC1		0.52	0.55	0.52	0.54	0.39	0.47	0.43
S 1	0.80		0.64	0.61	0.59	0.48	0.53	0.51
S 2	0.81	0.93		0.64	0.65	0.49	0.55	0.54
IP1	0.79	0.93	0.93		0.62	0.49	0.52	0.52
IP2	0.85	0.92	0.96	0.95		0.44	0.51	0.48
P1	0.68	0.82	0.81	0.83	0.76		0.52	0.55
P2	0.74	0.83	0.83	0.80	0.80	0.90		0.58
P3	0.68	0.80	0.83	0.81	0.77	0.97	0.94	
Statistic								
No. Items	33	22	31	21	29	18	21	25
Mean SS	8.43	8.25	8.23	8.14	8.15	8.52	8.37	8.61
SD SS	0.82	1.06	0.96	1.13	0.93	0.95	1.02	0.84
Mean Raw	26.38	17.05	23.96	15.98	22.10	14.60	16.60	20.70
SD Raw	3.54	3.08	3.93	3.12	3.57	2.29	2.85	2.79
Reliability	0.64	0.66	0.71	0.67	0.64	0.53	0.63	0.61

Table 6. RAD Section Score Correlation Matrix and Statistics

Section Name
Patient Care
Safety
Image Production
Procedures
Patient Interactions and Management
Radiation Physics and Radiobiology
Radiation Protection
Image Acquisition and Technical Evaluation
Equipment Operation and Quality Assurance
Head, Spine and Pelvis Procedures
Thorax and Abdomen Procedures
Extremity Procedures

RAD Section Key:



Content Area	PC1	S 1	IP1	P1	P2	P3	P4	P5
PC1		0.39	0.51	0.46	0.41	0.51	0.43	0.48
S 1	0.66		0.57	0.59	0.54	0.53	0.51	0.54
IP1	0.77	0.85		0.66	0.60	0.64	0.61	0.61
P1	0.75	0.95	0.96		0.58	0.61	0.57	0.57
P2	0.66	0.87	0.87	0.91		0.61	0.60	0.59
P3	0.79	0.81	0.89	0.91	0.90		0.63	0.60
P4	0.72	0.84	<i>0.91</i>	0.92	0.96	0.97		0.58
P5	0.76	0.85	0.87	0.88	0.90	0.86	0.90	
Statistic								
No. Items	20	22	38	24	24	28	20	24
Mean SS	8.31	8.08	8.10	8.38	8.55	8.44	8.34	8.48
SD SS	0.87	0.93	0.84	0.85	0.80	0.86	0.92	0.87
Mean Raw	14.91	15.58	27.10	18.13	18.78	21.41	15.01	18.53
SD Raw	2.70	3.13	4.95	3.20	3.02	3.74	2.88	3.29
Reliability	0.59	0.61	0.73	0.64	0.65	0.71	0.60	0.67

Table 7. NMT Section Score Correlation Matrix and Statistics

NMT Section Ke	y.
Abbreviation	Section Name
PC	Patient Care
S	Safety
IP	Image Production
Р	Procedures
PC1	Patient Interactions and Management
S 1	Radiation Physics, Radiobiology, and Regulations
IP1	Instrumentation
P1	Radionuclides and Radiopharmaceuticals
P2	Cardiac Procedures
P3	Endocrine and Oncology Procedures
P4	Gastrointestinal and Genitourinary Procedures
P5	Other Imaging Procedures

Content Area	PC1	PC2	S 1	S2	P1	P2	P3	P4
PC1		0.40	0.46	0.49	0.44	0.39	0.41	0.49
PC2	0.84		0.42	0.47	0.52	0.40	0.43	0.51
S 1	0.84	0.78		0.61	0.46	0.47	0.54	0.57
S 2	0.91	0.88	1.00		0.53	0.47	0.55	0.59
P1	0.83	1.00	0.77	0.90		0.49	0.47	0.59
P2	0.83	0.85	0.89	0.89	0.95		0.48	0.51
P3	<i>0.79</i>	0.84	0.92	0.95	0.83	0.96		0.57
P4	0.88	0.93	0.90	0.95	0.96	0.95	0.96	
Statistic								
No. Items	25	22	20	29	26	18	24	36
Mean SS	8.34	8.25	7.97	8.02	8.21	8.36	8.11	8.23
SD SS	0.84	0.90	1.14	0.95	0.90	0.94	0.95	0.82
Mean Raw	19.91	17.27	14.93	21.89	20.38	14.39	18.39	28.14
SD Raw	2.66	2.51	2.88	3.47	3.01	2.15	2.89	3.77
Reliability	0.48	0.47	0.62	0.60	0.58	0.46	0.55	0.64

Table 8. THR Section Score Correlation Matrix and Statistics

THR	Section	Key:
-----	---------	------

THIC Dection Re	<i>j</i> •				
Abbreviation	Section Name				
PC	Patient Care				
S	Safety				
Р	Procedures				
PC1	Patient Interactions				
PC2	Patient and Medical Record Management				
S 1	Radiation Physics, Equipment, and Quality Assurance				
S2	Radiation Protection				
P1	Treatment Sites and Tumors				
P2	Treatment Volume Localization				
P3	Prescription and Dose Calculation				
P4	Treatments				

Content Area	PC1	IP1	IP2	IP3	P1	P2	P3
PC1		0.42	0.46	0.49	0.41	0.46	0.44
IP1	0.57		0.79	0.66	0.62	0.65	0.58
IP2	0.64	0.94		0.67	0.59	0.63	0.56
IP3	0.76	0.87	0.90		0.68	0.65	0.65
P1	0.55	0.71	0.69	0.88		0.83	0.79
P2	0.61	0.73	0.72	0.82	0.91		0.77
P3	0.65	0.73	0.72	0.92	0.97	0.93	
Statistic							
No. Items	29	50	44	21	75	109	32
Mean SS	8.06	7.37	7.87	8.12	7.93	8.00	7.62
SD SS	0.91	1.24	1.09	1.15	1.01	1.07	1.10
Mean Raw	21.72	32.66	31.75	15.88	54.26	80.56	22.10
SD Raw	3.35	8.05	6.19	3.15	10.06	15.63	4.72
Reliability	0.62	0.86	0.82	0.68	0.89	0.93	0.74

Table 9. SON Section Score Correlation Matrix and Statistics

SON Section Ke	y:
Abbreviation	Section Name
PC	Patient Care
IP	Image Production
Р	Procedures
PC1	Patient Interactions and Management
IP1	Basic Principles of Ultrasound
IP2	Image Formation
IP3	Evaluation and Selection of Representative Images
P1	Abdomen
P2	OB/GYN
P3	Superficial Structures and Other Sonographic Procedures



Content	PC1	S 1	IP1	IP2	IP3	P1	P2	P3
Area		0.40	0.52	0.51	0.50	0.50	0.51	0.26
PC1		0.49	0.52	0.51	0.50	0.50	0.51	0.36
S 1	0.85		0.66	0.67	0.65	0.58	0.52	0.40
IP1	0.75	0.95		0.80	0.75	0.62	0.61	0.47
IP2	0.73	0.96	0.95		0.79	0.70	0.63	0.51
IP3	0.73	0.96	0.92	0.96		0.65	0.62	0.44
P1	0.76	0.88	0.78	0.87	0.83		0.68	0.57
P2	0.82	0.84	0.82	0.84	0.85	0.96		0.50
P3	0.71	0.78	0.77	0.82	0.73	0.97	0.91	
Statistic								
No. Items	17	15	40	38	34	26	20	10
Mean SS	7.72	8.03	8.09	8.15	7.84	8.33	7.84	8.22
SD SS	1.16	1.18	1.12	1.19	1.15	1.09	1.17	1.24
Mean Raw	11.88	11.06	29.93	28.72	24.30	20.30	14.31	7.67
SD Raw	2.67	2.37	6.06	6.02	5.30	3.82	3.17	1.69
Reliability	0.58	0.58	0.83	0.84	0.80	0.76	0.66	0.45

Table 10. MRI Section Score Correlation Matrix and Statistics

NMT Section Key:					
Abbreviation	Section Name				
PC	Patient Care				
S	Safety				
IP	Image Production				
Р	Procedures				
PC1	Patient Interactions and Management				
S 1	MRI Screening and Safety				
IP1	Physical Principles of Image Formation				
IP2	Sequence Parameters and Options				
IP3	Data Acquisition and Processing				
P1	Neuro				
P2	Body				
P3	Musculoskeletal				

When interpreting the correlations in Tables 6 through 10, it is important to consider the reliability of each section score. Sections with low reliability will have low correlations with other subscales. This is why the report provides the corrected correlations. A low reliability coefficient for a section also indicates that a candidate's score for that section is only an approximation of the candidate's true level of knowledge. For this reason, ARRT cautions students and program directors not to over-interpret small score differences among section scores. The limited reliability of section scores is the primary reason that ARRT bases its pass/fail decisions on total scores. Total scores are sufficiently reliable to make pass/fail decisions; section scores may not have sufficient reliability to make those decisions. A notable exception to this is SON. ARRT does base pass/fail decisions on the Abdomen and OB/GYN sections of that exam, and the reliability of those section scores is quite high.

Reliability of Exam Scores

Reliability refers to the repeatability and consistency of exam scores. A candidate who takes one form of an exam on one occasion and a second parallel form on another occasion should earn similar scores if the exam scores are reliable and the candidate has not changed in the time between the exam administrations (i.e., learned new material). Major differences should occur only if there is true change in the candidate's knowledge or if the exam scores are unreliable.

Reliability also describes how well candidates' observed scores on an exam approximate their "true" scores. A candidate's true score may be defined as the mean of their observed scores from a very large number of examinations. The true score is theoretical and not observable in practice.

Reliability coefficients are estimates of the reliability of exam scores. Reliability coefficients typically range from zero to one, with values near one indicating high consistency and those near zero indicating little or no consistency. In this report, Cronbach's coefficient α is the reliability estimate of choice. Cronbach's α , which requires only one exam administration, is an estimate of the reliability of a group's exam scores. Although it is never possible to determine the exact amount of error in one specific candidate's score, the standard error of measurement (SEM) describes the expected variation of each candidate's observed score around that candidate's true score.

Coefficient Alpha

The equation for Cronbach's coefficient α is

$$\alpha = \frac{k}{k-1} \left(1 - \frac{\sum_{i=1}^{l} \hat{\sigma}_{i}^{2}}{\hat{\sigma}_{X}^{2}} \right), \tag{1}$$

where *k* is the number of items,

I is the total number of items,

X is a set of exam scores,

 $\hat{\sigma}_i^2$ is the variance on an individual item *i*, and

 $\hat{\sigma}_x^2$ is the total exam variance.

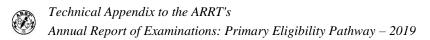


Table 11 contains the reliability estimates for RAD, NMT, THR, SON, and MRI. Recalling that reliability coefficients range from 0.0 to 1.0, one can see that the reliability estimates for the exam scores are quite high at 0.90 or greater. These high reliability estimates mean that observed scores for these exams likely correspond quite closely to true scores for these exams.

Tuble 11. Mean makes of meethal consistency and Standard Error of Measurement						
Discipline	a	SEM at the	e Mean Score	SEM at the	SEM at the Cut Score	
	α	Raw	Scaled	Raw	Scaled	
RAD	0.93	5.46	2.06	6.27	2.37	
NMT	0.93	5.73	1.85	6.43	2.07	
THR	0.90	5.57	2.20	6.22	2.45	
SON	0.97	7.98	1.65	8.38	1.74	
MRI	0.95	5.87	2.17	6.34	2.35	

Table 11. Mean Indices of Internal Consistency and Standard Error of Measurement

Standard Error of Measurement

The standard error of measurement (SEM) is a type of standard deviation. SEM is the standard deviation of a hypothetical set of repeated measurements for a single individual. A common equation calculates the SEM using the reliability estimate, r_{XX} (α from Equation 1), and the standard deviation of exam scores, S_X , with the equation

$$SEM = S_X \sqrt{1 - r_{XX}} .$$
 (2)

The above equation for SEM represents the mean SEM across all exam scores. SEM is not consistent, however, across the full range of scores, especially at the extremes. The SEM calculated at the cut score and the mean score will give a more accurate picture of the standard error. The equation for SEM at a particular score is

$$\text{SEM}_{\hat{X}} = \sqrt{\left(\frac{\hat{X}(k-\hat{X})}{k-1}\right)\left(\frac{1-r_{XX}}{1-r_{21}}\right)},$$
(3)

where \hat{X} is a score value of interest,

k is the number of items,

 r_{XX} is the reliability of scores using Cronbach's α , and

r₂₁ is the reliability of scores using Kuder-Richardson Equation 21 (Lord, 1955; Keats, 1957).

Table 11 provides the standard error of measurement for the mean score and the cut score in both raw and scaled score units using Equation 3.



Decision Consistency

ARRT administers examinations with criterion-referenced cut score standards as the basis of decisions to grant certification and registration. Agreement indices quantify the consistency or reproducibility of those dichotomous (two option) decisions. Decision consistency in this case describes how consistently the examinations classify individuals into certified and registered and not certified and registered groups. When organizations base a pass/fail decision on a single exam score, there will be a small number of candidates who passed but should have failed (false positives) and a small number of candidates who failed but should have passed (false negatives). The threshold loss agreement indices used in this report focus on the consistency of classifications, treating all potential misclassification errors as equally serious.

The threshold loss indices assume a dichotomous, qualitative classification of candidates as certified and registered or not certified and registered based on a cut score. The methods were originally developed using two or more exam administrations for every candidate. Because multiple examinations are not practical, researchers developed alternative methods to estimate the indices with a single exam administration. This report uses a method developed by Subkoviak (1976) to estimate two threshold loss indices, p_0 and kappa. The estimation procedure assumes that a candidate's observed scores are independently and binomially distributed according to the number of exam items and the candidate's proportion-correct true score.

p_0 index

The p_0 index measures the overall consistency of pass/fail classifications. It is the proportion of individuals expected to be consistently classified as certified and registered and not certified and registered based on Subkoviak's (1976) method. The index is sensitive to the cut score, exam length, and score variability. For example, p_0 values will be smaller for cut scores near the mean of scores, because there are more people located near the mean than at the extremes if scores are normally distributed. The first column in Table 12 contains the p_0 values for each of the exams that this report covers. Classification decisions based on these exams are consistent between 89% and 94% of the time. This is a high level of decision consistency.

Table 12. Threshold Loss Indices					
Discipline	p_0	p_c	kappa		
RAD	0.94	0.80	0.70		
NMT	0.94	0.81	0.68		
THR	0.93	0.77	0.70		
SON*	0.89	0.60	0.73		
MRI	0.93	0.64	0.81		

* The p_0 statistic for SON makes a statistical adjustment to Subkoviak's (1976) method that takes into account the necessity to pass the overall exam, the Abdomen section, and the OB/GYN section.



Карра

While high classification consistencies are good, it is possible that some or many of the correct classifications of certified and registered or not certified and registered were due to chance. For example, a person can correctly guess heads or tails at the flip of a coin a certain percentage of the time. These correct guesses are due purely to chance. Kappa is a statistical index that shows the proportion of individuals consistently classified beyond that expected by chance. The equation for kappa is

$$k = \frac{p_0 - p_c}{1 - p_c},\tag{4}$$

where p_0 is the overall consistency of certified and registered/not certified and registered classifications and p_c is the proportion of consistent classifications that would be expected by chance.

The calculation for p_c is simply

$$p_c = (P_{Pass})^2 + (1 - P_{Pass})^2,$$
(5)

where P_{pass} is the proportion of people who pass the exam (Croker & Algina, 1986). Table 10 contains the kappa statistics for ARRT's exams. The kappa coefficient indicates that ARRT's exams consistently classify between 68% and 81% of the candidates above and beyond those already correctly classified by chance.

With regard to psychometric properties, ARRT's examinations are comparable to other welldeveloped examinations. ARRT's exam scores are reliable, with α coefficients at or above .90. The threshold loss indices indicate that most candidates are consistently classified as either certified and registered or not certified and registered. Maintaining a high-quality examination program is a vital part of ARRT's mission of promoting high standards of patient care by recognizing qualified individuals in medical imaging, interventional procedures, and radiation therapy. The results from this technical report show that ARRT indeed continues to develop quality examinations.

References

- Crocker, L., & Algina, J. (1986). *Introduction to classical and modern test theory*. Belmont, CA: Wadsworth.
- Keats, J.A. (1957). Estimation of error variances of test scores. Psychometrika, 2, 29-41.
- Lord, F.M. (1955). Estimating test reliability. *Educational and Psychological Measurement*, 15, 325-336.

Nunnally, J.C. (1978). Psychometric theory. New York: McGraw-Hill.

Subkoviak, M.J. (1976). Estimating reliability from a single administration of a mastery test. *Journal of Educational Measurement, 13*, 265-276.