

Data Sets for the Qualification of Volumetric CT as a Quantitative Imaging Biomarker in Lung Cancer[◇]

March 2010

Andrew J. Buckler, M.S., Buckler Biomedical LLC, and Chair, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Lawrence Schwartz, M.D., Chair, Department of Radiology, Columbia University, and Co-Chair, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Nicholas Petrick, Ph.D., Deputy Director, Center for Devices and Radiological Health, Food and Drug Administration, and Member, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Michael McNitt-Gray, Ph.D., DABR, Professor of Radiological Sciences, David Geffen School of Medicine at UCLA and Member, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Binsheng Zhao, DSc, Associate Professor and Director, Computational Image Analysis Lab, Department of Radiology, Columbia University, and Member, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Charles Fenimore, Ph.D., Information Technology Laboratory, National Institute of Standards and Technology, and Member, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Anthony P. Reeves, Ph.D., Professor, Electrical and Computer Engineering, Cornell University, and Member, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

P. David Mozley, M.D., Merck Research Laboratories and Chair, Extended Pharma Imaging Group and Co-Chair, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

Ricardo S. Avila, M.S., Senior Director of Healthcare Solutions, Kitware Inc., and Member, Quantitative CT Technical Committee, Quantitative Imaging Biomarker Alliance

[◇]**Datasets associated with this article are available at**

<http://midas.osa.org/midaspre/midas/item/view/912?key=a3dHdkVmOHpJWjlyTQ==>



Abstract

The drug development industry is faced with increasing costs and decreasing success rates. New ways to understand biology as well as the increasing interest in personalized treatments for smaller patient segments requires new capabilities for the rapid assessment of treatment responses. Deployment of qualified imaging biomarkers lags apparent technology capabilities. The lack of consensus methods and qualification evidence needed for large-scale multi-center trials, as well as the standardization that allows them, are widely acknowledged to be the limiting factors. The current fragmentation in imaging vendor offerings, coupled with the independent activities of individual biopharmaceutical companies and their contract research organizations (CROs), may stand in the way of the greater opportunity were these efforts to be drawn together. A preliminary report of the Quantitative Imaging Biomarkers Alliance (QIBA) activity was presented at a meeting of the Extended PhRMA Imaging Group sponsored by the Drug Information Agency (DIA) in October 2008.¹ The clinical context in Lung Cancer and a methodology for approaching the qualification of volumetric CT as a biomarker has since been reported.^{2,3} This report reviews the effort to collect and utilize publicly available data sets to provide a transparent environment in which to pursue the qualification activities in such a way as to allow independent peer review and verification of results. This article focuses specifically on our role as stewards of image sets for developing new tools.

Key words: quantitative imaging; therapy response; imaging biomarker; volumetric CT; regulatory pathway

Unmet Medical Needs as Business Drivers for Qualifying Quantitative Imaging

Problems with qualitative impressions of longitudinal changes in tumor burden before and after treatment include inadequate levels of inter-reader concordance when responses are less than dramatic. Discordance among "readers" has led to skepticism about medical imaging as a biomarker of response, as well as confusion about whether some investigational new drugs should be approved for general use.

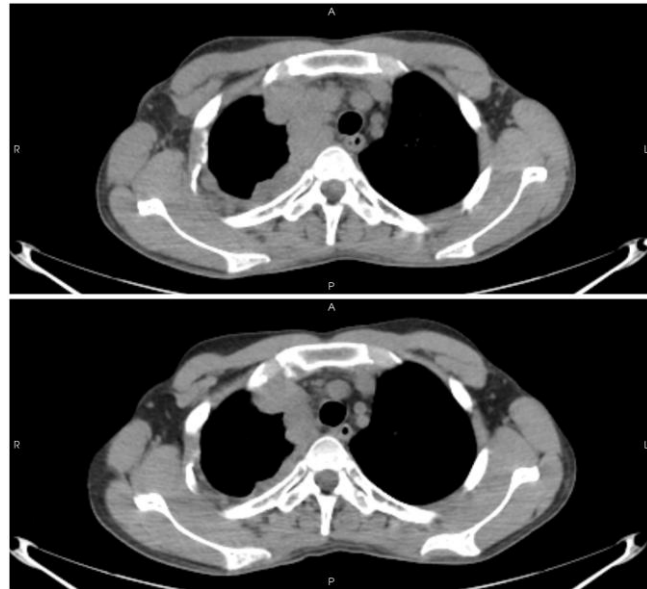
Subjective impressions are sufficient when the impact of treatment is so robustly effective that changes are conspicuous, just as they are when the therapy fails so completely that disease progression is obvious. However, as the "war on cancer" matures from our initial hopes of curing the disease into aspirations for converting patient management from acute therapy to manage morbidity over progressively longer and longer time horizons, needs for rapidly assessing the incremental value of adding new drugs to the standard of care are becoming increasingly important.

For an individual patient in an ordinary medical setting, being prescribed a marketed treatment regimen that has been established as sufficiently safe and effective in large populations is analogous to starting a personal clinical trial. This is because even the best treatment regimens fail in a some portion of patients with the disease, and even relatively safe therapies cause serious side effects in some people. These principles seem to hold for all treatments, and particularly for anti-neoplastic therapies. Patients want to know as soon as possible if their new-to-them treatment is conveying benefits. If it is not, then they want to launch a search for alternatives as soon as possible.

No one wants to waste time, effort, and money on treatments that are not helpful. From this perspective, the interests of individual patients and third party payers seem highly concordant. Many new treatments are expensive. Some are cost effective in individuals, but less so in large populations. New methods are needed to determine who is who. Until definitive enrichment tools are developed for matching individual patients to specific treatments, the early assessment of response will remain the primary mechanism for sparing resources.

Biopharmaceutical enterprises view clinical trials of novel products the same way as the other stakeholders in the management of cancer. Like individual patients, industry wants its products to succeed for the patients who use them, and as a consequence, produce a net-positive return on investment. More sensitive biomarkers of response would allow industry to reduce the number of patients required to test new products, as well as decrease the amount of time that patients need to remain on-study. The net effect would increase the number of new treatments for unmet medical needs that reach the market and make a positive impact on human health, primarily by allowing investigational new treatments to fail faster than is currently possible in clinical trials that use survival or clinical signs of progression as their endpoints.

Spatially specific biomarkers could provide more informative data than clinical outcomes in patients with heterogenous metastases. Consider the case shown in Figure 1. This 21 year old man presented with a chief complaint of shortness of breath. Panel 1a shows a mass compressing the right lung and displacing the trachea to the left. Panel 1b shows the beneficial effects of monotherapy with an experimental agent. After 18 days, tumor volume decreased, the trachea moved back towards the midline, and the patient reported symptomatic relief of dyspnea.



(a)

(b)

Figure 1: A lung cancer patient's CT scans before (a) and after (b) the administration of an experimental monotherapy ([View 1](#)).

89 Figure 2 shows that some, but not all, tumors in the chest became larger and more metabolically active
 90 at the same time others moved towards remission. In fact, this patient came off trial after only 6 weeks
 91 because a new metastasis caused a spinal cord compression despite the fact that the masses causing
 92 dyspnea continued to show a favorable response. In this case, relying on clinical outcomes alone would
 93 have led to a conclusion that the drug is not active because the patient failed treatment after only 6
 94 weeks. But, a more scientifically accurate conclusion might be that the drug holds promise for treating
 95 some tumor populations, but not others.

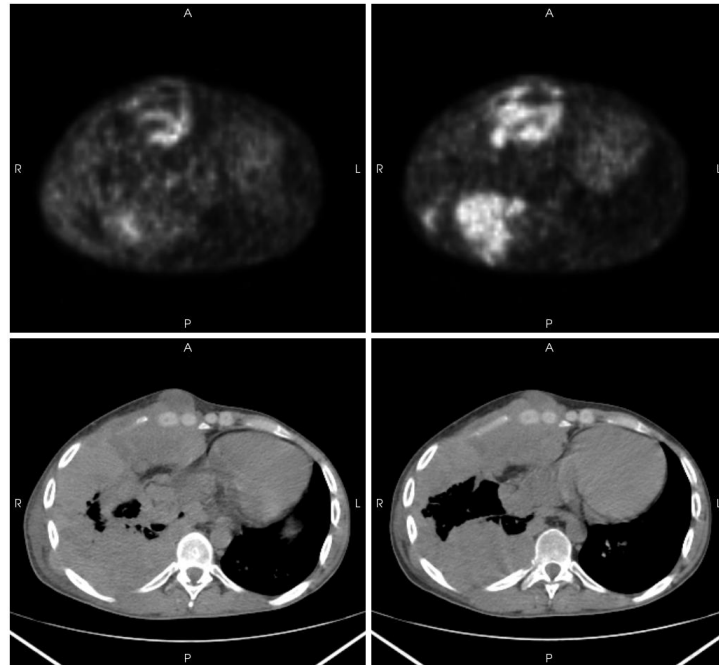


Figure 2: PET (top) and CT (bottom) scans of the lung cancer patient shown in Figure 1 at baseline (left) and 18 days after (right) the administration of therapy ([View 2](#)).

If we imagine a future state where quantitative imaging provides regionally specific information about tumor responses in the whole patient that triggers the addition of treatment options to therapies that are providing selective benefits to some tumors but not others, then the imperative for qualifying quantitative imaging biomarkers becomes easier to visualize.

Response Evaluation Criteria in Solid Tumors (RECIST) is a quantitative image analysis technique based on diameter measurements selected from axial slices. It is designed to meet these needs, particularly when responses are robust. However, problems with diameter measurements on axial slices include their lack of sensitivity. The categorical response of Stable Disease is so broad that classifying a treatment as effective or futile can take a long time. This is in part because thresholds for categorical responses correspond to changes in volume of about -66% for partial response to about 73% for progressive disease. While changes in longest diameters have never been validated or qualified in the formal sense we are pursuing for volumetrics, RECIST has been used in a very broad number of cases and is generally recognized as effective for tumors that tend to have spherical geometries and contract or expand more or less uniformly. Unfortunately, the more complex the tumor morphology and pattern of longitudinal change, the less sensitive the formalism becomes, to the point where it can be misleading in some cases.

Referring to figure 3, changes in the longest diameters of the target lesions suggested that this patient remained in a prolonged state of Stable Disease. As a consequence, the subject added little analytical power needed to distinguish between the two arms of the trial. In retrospect, volumetric image analysis suggests that this patient had an initial response to treatment, but could have come off trial and switched to a new treatment several months before changes in unidimensional line-lengths met criteria for Progressive Disease.

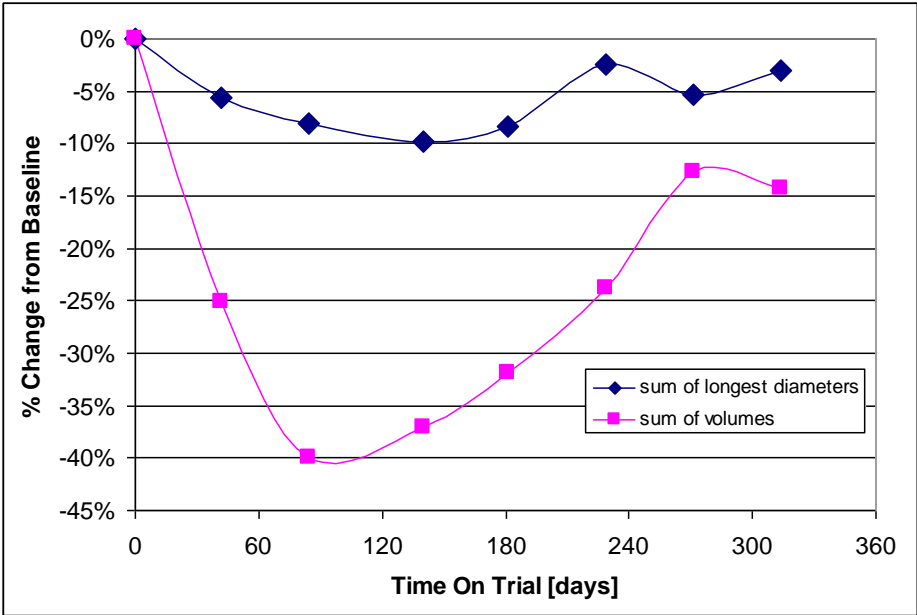


Figure 3

All of the stakeholders lose when benefits are not recognized, or there is a delay in diagnosing Progressive Disease.

Methods

It is widely recognized that significant advances in imaging technology have led to an increasingly important role for imaging in diagnosis, staging, guiding systemic, local, or interventional therapies, and monitoring responses to treatment. However, development of imaging technologies is expensive, and early phase justification of effectiveness, before commercial viability is established, can be difficult. There is an emerging consensus that a cooperative atmosphere must be developed among the biopharmaceutical industry, the imaging device manufacturers, government funding agencies, and regulatory authorities, as well as scientists in a wide range of fields, to cost effectively select and qualify mature quantitative imaging methods as biomarkers for the measurement of response to therapy.

The development of public resources and open source tools for imaging as a biomarker using X-ray CT was re-invigorated by the NCI, NIBIB, FDA and National Institute of Standards and Technology (NIST) in 2005, which included collaboration with the Radiological Society of North America (RSNA).^{4,5,6,7,12} This earlier work prompted the organization of an inter-federal agency workshop held at NIST in September 2006, which addressed physical standards for imaging as a biomarker.² Stakeholders from academia, industry, and scientific imaging societies including RSNA, American Association of Physicists in Medicine (AAPM), Society of Nuclear Medicine (SNM), and International Society for Magnetic Resonance in

Medicine (ISMRM) proposed a model similar to the “Integrating the Healthcare Enterprise” (IHE) paradigm to engage industry stakeholders in this research area.

At its annual meeting in 2007, RSNA created the Quantitative Imaging Biomarker Alliance (QIBA) to investigate the role of quantitative imaging methods in CT, MRI and PET as potential biomarkers in evaluating disease and responses to treatment. The alliance has formed technical committees of representatives from the instrumentation manufacturers, software developers, imaging professionals in the pharmaceutical industry, radiologists from the imaging contract research organizations (CROs), officers in regulatory agencies, governmental research organizations, imaging scientists, and professional imaging society representatives. One of the technical committees is referred to as the “Quantitative CT Technical Committee.”

The Quantitative CT Technical Committee is engaged to produce alternative methods of response assessment, based on volumetric image acquisition and analysis, which will be accepted through appropriate regulatory pathways as predictors of clinical benefits, such as overall survival (OS). The first specific aim compares time-dependent outcome measures based on uni-dimensional longest diameters to analogous endpoints based on 3D volumetric image analyses. The expectation is that these alternative methods would be adopted if they require fewer enrollees in clinical trials, shorten time on trial for each subject who will ultimately fail to benefit from treatment, decrease the length of time required to conduct trials, and/or provide better correlations with actual clinical outcomes.

The Committee was formed to include practicing clinicians, professional society leaders, regulatory officers, pharmaceutical industry representatives, imaging scientists, and imaging device industry representatives. The principal value of the effort is to help converge the interests and effort of many stakeholders.

Long-Term Goals are to establish processes and profiles that will eventually lead to the acceptance by the imaging community, clinical trial industry, and regulatory agencies, of 3D volumetric CT as *proof of biology, proof of changes in pathophysiology, and surrogate end-points for changes in the health status of patients.*

Specific Aims are to develop the capability to meet targeted levels of accuracy and reproducibility for the quantification of anatomical structures, such as neoplastic masses. This in turn requires identifying and creating coping strategies for all significant sources of variability in these measurements.

Context is that this work is being conducted under the aegis of the RSNA's QIBA in collaboration with FDA's Division of Applied Math/ Office of Science and Engineering Laboratories (OSEL)/ Center for Devices and Radiological Health (CDRH), NCI, NIST, American College of Radiology Imaging Network (ACRIN), major imaging equipment manufacturers (Philips, GE, Siemens, Toshiba, etc.), the Extended Pharmaceutical Research and Manufacturers of America (PhRMA) Imaging Group, and others.

Constraint is that this work depends on the collaboration of, and must demonstrate benefit to, the imaging industry, the pharma industry, the academic research community, individuals with cancer, and the clinical community. The benefits must be robust to justify the increased time and effort required when compared to qualitative impressions, as well as satisfy the requirements of the regulatory agencies. Our approach is to converge scientific analysis in a way that encourages vendor participation while meeting current biopharmaceutical industry needs.

Our ultimate goal is the use of these biomarkers on typical imaging systems in the practice of medicine.

Results to Date

The QIBA initiative has explored a number of issues and opportunities to improve research and development of volumetric CT therapy assessment methods. To accomplish this, it has been essential to obtain and analyze a wide range of image data collections that span clinical concepts and challenges, fundamentals of image acquisition, and opportunities to better perform the evaluation of algorithm performance. The sections that follow describe these data collections and the important insights each collection provides to the research community.

Understanding Performance on Phantoms

One approach to efficiently develop and evaluate the applicability of a quantitative imaging biomarker is to investigate the biomarker's performance with phantom data. Phantom image data can come in many forms including imaging simple lesion-like objects on flat backgrounds or imaging anthropomorphic phantoms containing realistic structure, complex synthetic lesions, and realistic physiology. Figure 4 shows three different examples of lung and chest phantoms from the literature, including a tissue equivalent tissue equivalent thorax section phantom (Fig. 4a), an anthropomorphic chest phantom (Figure 4b, and a mechanical breathing phantom (Figure 4c).^{8,9,10}

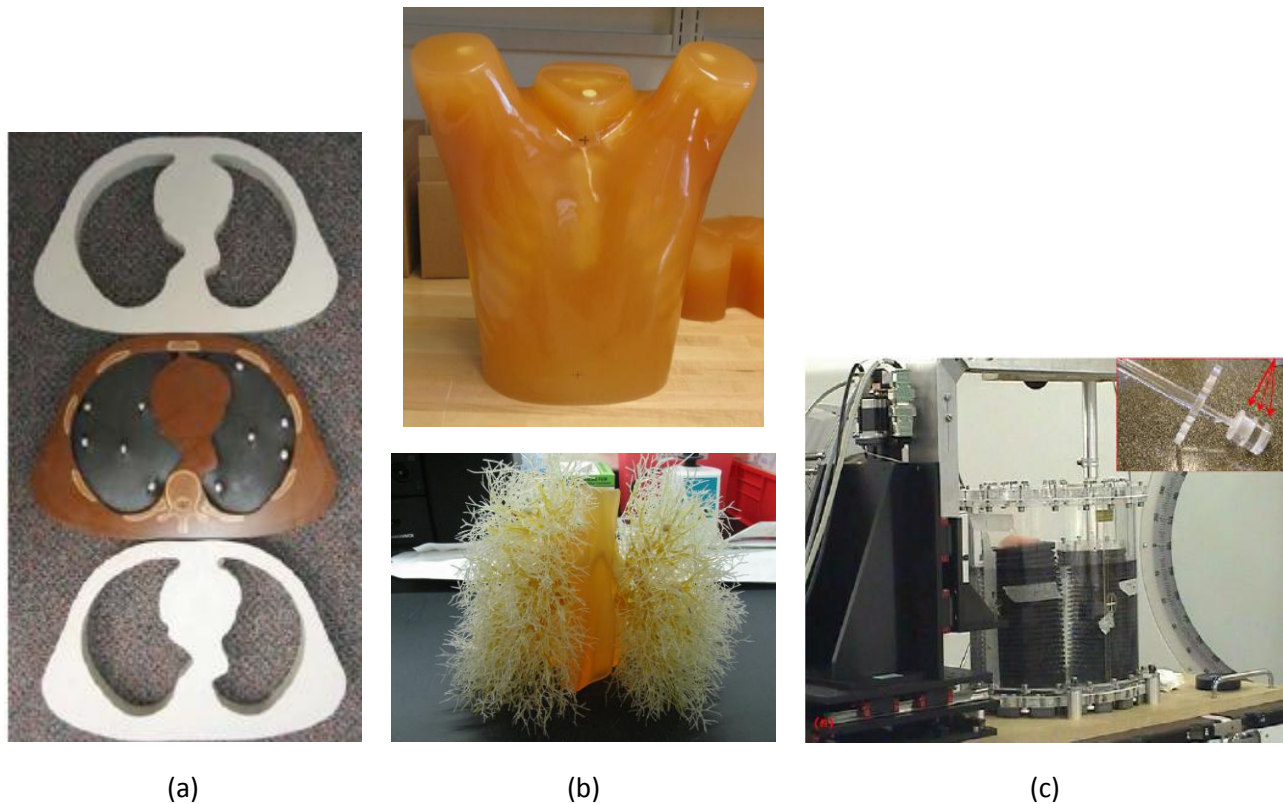


Figure 4: (a) tissue equivalent thorax section phantom (center) containing 9.5 mm diameter simulated spherical lung nodules, with two water-equivalent bolus sections (top and bottom), (b) the exterior shell of an anthropomorphic thoracic phantom and its vasculature insert; and (c) a mechanical lung phantom used to simulate breathing. Images in (a)-(c) are reprinted with permission from Refs. 8-10, respectively.

Although phantoms are different from real patients in many ways, phantom studies allow for a systematic analysis of biomarker performance against a known reference standard and under a range of

196 imaging conditions. This type of systematic analysis would be virtually impossible to conduct using
197 patient scans because of dose concerns, variability in patients, motion artifacts, and lack of a definitive
198 truth standard.¹¹ While phantom studies are unlikely to serve as a complete replacement for evaluating
199 a new biomarker on patient data, they may serve at least three important functions. One is to quickly
200 triage potential imaging biomarkers, so that time is not wasted evaluating biomarkers that have little
201 potential for providing reliable quantitative measurements. New biomarkers that don't perform well
202 with idealized phantom data are unlikely to perform well in patients whose diseases are well modeled
203 by the phantom. For those imaging biomarkers that do show promise, a second function of phantom
204 data could be to systematically probe how biomarker performance is impacted by variations in imaging
205 hardware and image acquisition protocols. Again, this type of systematic evaluation of a biomarker is
206 virtually impossible to conduct with patient data, even within a clinical trial, because of the large
207 variability in manifestations of disease both within and among patients. Finally, a third contribution of
208 phantom studies could be in the design of clinical trials incorporating an imaging biomarker. By first
209 understanding how variations in image acquisition affect the reliability of the quantitative measurement
210 through phantom studies¹², it becomes possible to develop appropriate imaging standards as well as
211 determining a minimum number of patients required to overcome the variability implicit when
212 implementing the imaging biomarker. Additional patients, above this minimum level, would be
213 necessary to overcome patient variability as well as other sources of error in any particular trial.

214 A companion manuscript in this issue by Gavrielides et al. describes CT image data for an
215 anthropomorphic thorax phantom containing synthetic lung nodules.⁹ These data were collected by the
216 U.S. Food and Drug Administration (FDA) to evaluate various lesions size measurement algorithms, and
217 to develop a more complete understanding of how algorithm performance changes with variations in CT
218 acquisition protocols and imaging hardware. Figure 4(b) shows the thorax phantom and vasculature
219 lung inserts to which synthetic nodules were attached and then imaged within the dataset. The
220 phantom was scanned with a Philips 16-row scanner (Mx8000 IDT, Philips Healthcare, Andover, MA) and
221 a Siemens 64-row scanner (Somatom 64, Siemens Medical Solutions USA, Inc., Malvern, PA). The data
222 were collected using a factorial design so that a large number of combinations of exposure, pitch, slice
223 collimation, reconstruction kernels and slice thickness were collected for both simple spherical nodules
224 as well as more complex ovoid, lobulated and spiculated synthetic nodules. Figure 5 shows a complete
225 CT scan of the phantom with seven spherical nodules of various sizes and densities attached to the
226 vasculature insert.

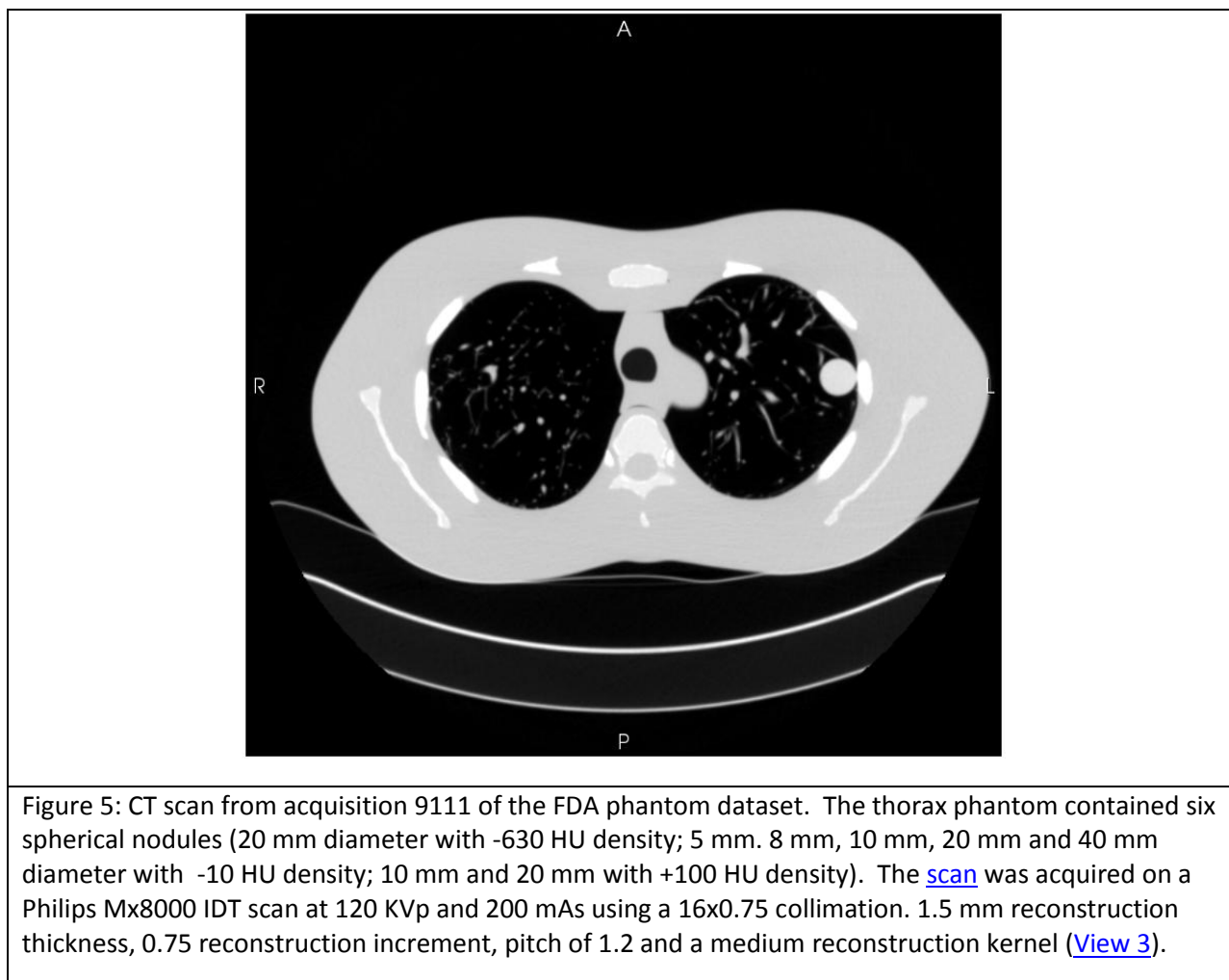


Figure 5: CT scan from acquisition 9111 of the FDA phantom dataset. The thorax phantom contained six spherical nodules (20 mm diameter with -630 HU density; 5 mm, 8 mm, 10 mm, 20 mm and 40 mm diameter with -10 HU density; 10 mm and 20 mm with +100 HU density). The [scan](#) was acquired on a Philips Mx8000 IDT scan at 120 KVp and 200 mAs using a 16x0.75 collimation. 1.5 mm reconstruction thickness, 0.75 reconstruction increment, pitch of 1.2 and a medium reconstruction kernel ([View 3](#)).

227

228 The FDA thorax phantom CT data described in Ref. 9 can be used as a resource for the development and
 229 assessment of lung nodule sizing algorithms. Both the bias and variance associated with a nodule sizing
 230 method can be obtained because the reference standard for nodule size as well as repeat exposures are
 231 included as part of the dataset. This makes the data ideal for comparing various size estimation
 232 algorithms. The data are also useful for developing new size estimation methods¹³ as well as developing
 233 appropriate assessment methodologies for comparing algorithms. These as well as various other
 234 applications of the phantom data are discussed in more detail in Ref. 9.

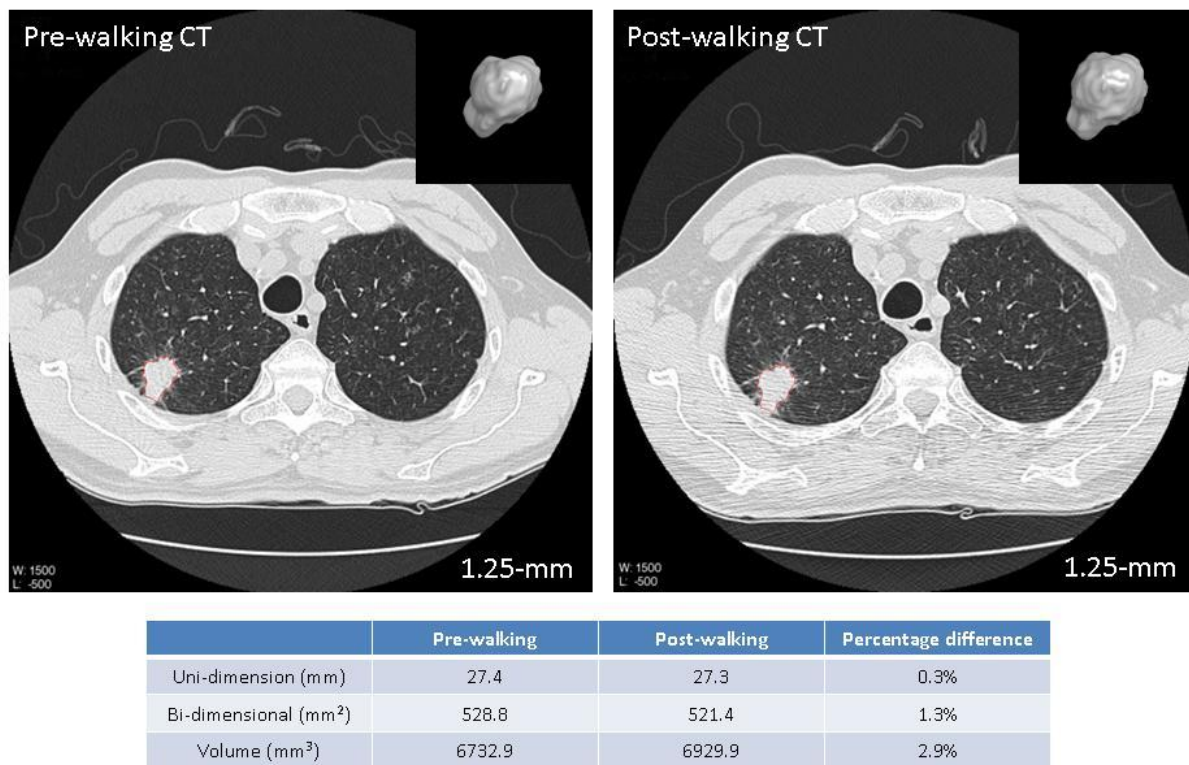
235 Evaluation of imaging biomarkers with phantom data is one important component in the qualification of
 236 these biomarkers in both drug trials and clinical practice. Clearly, phantom data have limitations
 237 because they do not match the diversity or complexity of real patients. This strongly suggests that
 238 testing on patient data will be necessary at some point in the development process, but also that
 239 phantom data can be a very effective tool in both streamlining the development process and maximizing
 240 the utility of patient image data.

241 **Clinical Data Resources**

242 There have been considerable efforts to create publicly available sets of image data to assist in some of
243 the efforts related to quantitative imaging of disease. These datasets represent an important aspect in
244 establishing quantitative imaging methods as they serve as reference datasets against which
245 investigators and researchers may be able to benchmark and compare their measurement algorithms.
246 Several datasets are now available, primarily through the NCI-funded Reference Image Database to
247 Evaluate Response to Therapy (RIDER).^{6,14,15}

248 Same-day repeat CT study in NSCLC patients

249 The first dataset to describe is the No-Change dataset provided by Memorial Sloan Kettering Cancer
250 center.¹⁶ In this study, 32 patients with Non-Small Cell Lung Cancer (NSCLC) were consented and
251 scanned twice within 15 minutes on the same scanner with the same imaging acquisition protocol. The
252 scanners were either LightSpeed 16 or VCT 64 (GE Healthcare, Milwaukee, WI). Images of each scan
253 were reconstructed at 1.25mm slice interval without overlap. This unique experiment represents repeat
254 scans under a presumed “no change” condition. Tumor differences measured between the two scans
255 can be considered as measurement variation/error that is possibly caused by intrinsic variance in the CT
256 scanning device, errors in the image processing system, differences in patient positioning, patient
257 inspiration level, etc. Because this dataset does contain the same lesions acquired on two repeat CT
258 scans under identical parameter settings in a short time period, it can be used to investigate minimum
259 detectable changes on the state-of-the-art CT scanners by using advanced measurement tools, the
260 information needed to define tumor response and progression. These datasets have been made publicly
261 available through the NBIA web archive (<http://ncia.nci.nih.gov/>).



Courtesy of Laboratory for Computational Image Analysis, Columbia University Medical Center

Figure 6: An example taken from the same-day repeat CT study. Computer-aided tumor measurements were different on the two repeat CT scans even if there were no biological change of the tumor ([View 4](#)).

CT lung studies at different time intervals

In another RIDER project related study, serial CT scan images of patients with known tumors in the lungs (both primary and metastatic lesions) were submitted to NBIA under the RIDER collection. Each case had at least 2 image data sets from different time points; many had 3 or more time points. These cases were collected from UT-MD Anderson Cancer Center and Memorial Sloan-Kettering Cancer Center, as part of their clinical operation. There was no specific attempt to tightly control the imaging parameters between studies for these patients.

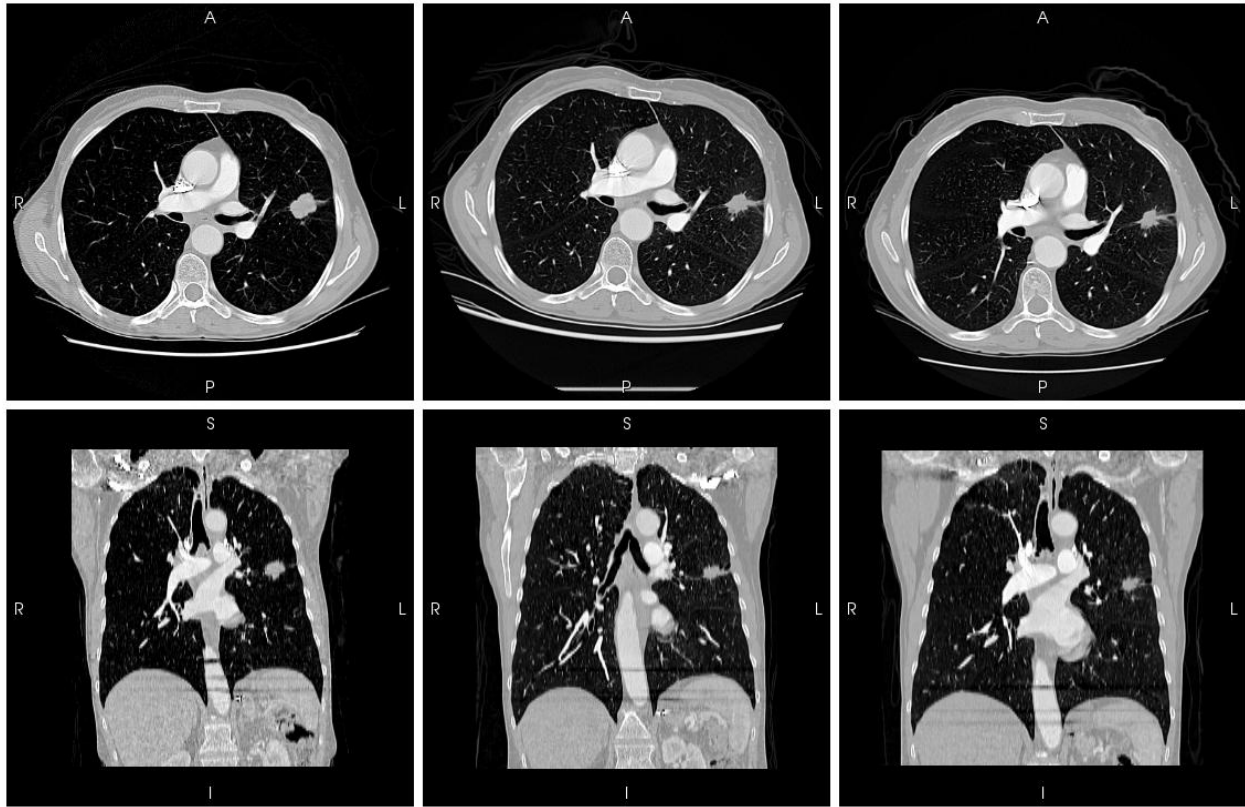


Figure 7: Longitudinal Scans where Patient has Known Tumor ([View 5](#)).

Another public resource for clinical CT image data is the Public Lung Database to Address Drug Response¹⁷. This dataset contains a number of different exemplar CT image sets including cases with at least two scans having manual volumetric boundary markings and cases with at least two scans recorded in the same session (zero-change) as part of a biopsy procedure that are documented with a semi-automated lesion measuring algorithm. These cases were collected from the Weill Cornell Medical College as part of their clinical operation.

While these reference datasets cannot be used to quantify the accuracy of measurement, they are a tremendous resource for researchers who need to characterize the precision of new quantitative imaging methods. They can be used to investigate the minimum detectable change (using the cases with no change) as well as different sources of variance (both sets).

Algorithm Evaluation Systems

We expect that computer assisted methods for measurement will aid the physician with respect to accuracy and precision of lesion measurements. One principal goal in evaluating such methods is to support the improvement of algorithms by providing developers a resource for identifying the strengths and weaknesses of their methods. Similar evaluations have been applied to computer vision methods for biometric-based identification, such as face and gait recognition.

For the clinical use of the volumetric image biomarker the most relevant measurement is the relative change in lesion size over some time interval. As has been stated before, it is critical to know when a measured change in size is statistically significantly greater than the measurement error (i.e., represents an actual change in the lesion); secondly we would like to know the precision of the size change

measurement. To explore these issues in the context of computer algorithms and real lesions rather than phantoms, studies have been conducted on selected data sets of pairs of lesions to evaluate how different computer algorithms compare on a standardized dataset.

A first evaluation of this type was Biochange'08, which invited participants to measure the change in pulmonary lesions using CT data from both the RIDER database of patients with known lung tumors and CT imaging of the FDA's anthropomorphic phantom described earlier.¹⁸ This pilot study provided algorithm and software developers with 13 cases, each having scans at 2 time points. Seven cases were clinical, all with 5.0 mm slice thickness and acquired at intervals of weeks to months. There were six phantom nodules from studies of the FDA phantom, having slice thicknesses of 3.0 mm and 0.8 mm.

The clinical data was chosen from 23 cases in RIDER for which diameter measurements on axial slices (one-dimensional) markup by 2 radiologists was available. In the analysis the markup was used as a reference and also examined the statistical differences between the algorithms/software.

The study was designed as a pilot, a proof of concept for the evaluation process. There were 3 participants who provided 4 submissions. Three of the submissions involved a software-assisted user in the loop. The study required the participant to submit a measure of change for each case. While this permitted the use of any change metric, for example one based on one- or two-dimensional measurement, each participant submitted the fractional change in volume and also provided volume measurements at both time points. The limited size of the study did not support statistically significant findings about the algorithms but did suggest some tentative conclusions regarding the comparison of diameter measurements on axial slices markup and computer assisted change measurement. The phantom data provided insight into the effects of slice thickness on the measurement of volume change.

The data suggests that the algorithms achieve agreement comparable to that between the radiologists and the two reach similar categorical conclusions. In particular, there were 6 cases for which the two radiologists agreed on the diameter measurements on axial slices categorical assessment (response/stable disease/disease progression) while, in one case, the radiologists disagreed. Using categorical (3-dimensional) thresholds derived from the diameter measurements on axial slices criteria, the 4 submissions obtained results similar to those of the radiologists: agreeing with each other and with the radiologists in 5 of the 6 cases. The two cases of disagreement occurred on lesions involved, in one case, with the mediastinum and in the other, with the lung wall at the apex. Figure 8 shows CT slices of the involved lesion near the lung apex on which the computed results disagreed. In this case the radiologists' markup agreed in finding stable disease.

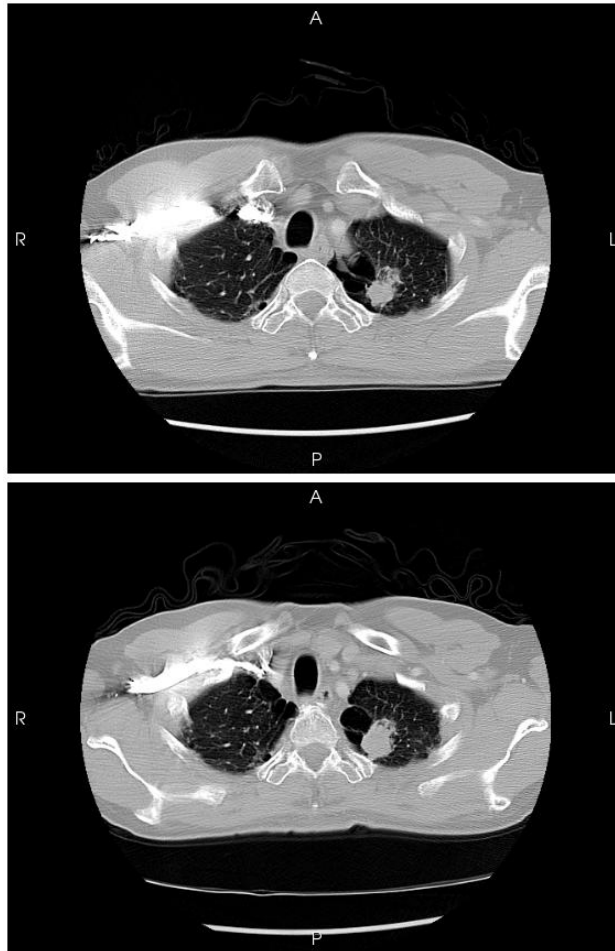


Figure 8: An involved apical lesion at two time points, 7 months apart. In this RIDER case, used in Biochange'08, four computer-assisted measurements did not agree on a categorical assessment of volume change akin to diameter measurements on axial slices. In mark up by two radiologists, the diameter measurements on axial slices criteria indicate stable disease. The four computer-assisted results agreed on categorical volume change for six other clinical cases ([View 6](#)).

The phantom nodules were scanned in both thin- and thick-slice series (0.8 and 3.0 mm). For the phantoms, there was no change. There was a striking difference between the thin and thick slice results. For thin slice, the absolute range of reported change measurements was less than 10%. For the thick slice data, the range was about 40%.

A follow-on study to the Biochange '08 pilot is the planned full scale Biochange Challenge. It also uses the RIDER lung CT studies but mainly has thin slice studies, including the MSKCC Coffee Break data discussed earlier. In addition to the participation of algorithm/software developers, the planned study seeks the participation of radiologists to provide markup for comparison with the computed change measures.

A second study group members have conducted is the "VOLCANO'09 Challenge."¹⁹ This challenge invited participants to evaluate the change in size of pulmonary nodules. The challenge involved measuring the change in nodule size for 50 scan pairs. Four additional scan pairs were made available for training. The data was selected from cases prepared for the Public Lung Database to Address Drug

Response.^{20,21} This database was sponsored by the Prevent Cancer Foundation²² and provides information on a number of aspects of lesion measuring by means of sample image; this resource is complimentary to the RIDER database. A key component of this database are repeat scans made at the same time. This zero change dataset is similar to the No-Change dataset except that scans were obtained from the start of CT guided biopsy procedure before the needle affects the image quality.

Teams reported the fractional change in nodule size for each of the 50 scan pairs. Thirteen different teams submitted their measurement change results from a total of 17 different methods. In 11 of these cases, the actual volumes recorded for each nodule were also reported. The participants were only informed that there were 50 nodule pairs; however, the data may be divided into four subgroups:

- A. (14) zero-change in which the scans were taken minutes apart and therefore there is no real change in the nodule size.
- B. (13) zero-change cases as in A above except that one scan has a slice thickness of 1.25 mm and the second scan has a larger slice thickness (2.5 or 5.0 mm)
- C. (19) nodules with a significant time interval between scans and therefore some real change and (3) nodules with a large amount of size change (greater than 1.5 times in volume). Of these nodules 19 were considered to be stable or benign by biopsy and 3 were diagnosed as malignant.
- D. (1) synthetic phantom nodule with a known size recorded with a different slice thickness+

If we only used zero-change data then any system that had a constant output set to zero would be considered to have an ideal response. For this reason we included cases for which a real change was indicated by observation; however, for these cases there is no way to know precisely how much that change is. Most evaluation methods for CAD systems, including challenges, involve a ground truth established by experts. However, for the task of nodule size estimation it is well known that there is a large amount of variation or disagreement in expert size estimations.²³ Further, it has not been established that expert's manual estimations are superior to automated measurements. In this challenge, while the change in size of nodules was reviewed by experts, the issue of ground truth was explored through the submitted responses to the challenge.

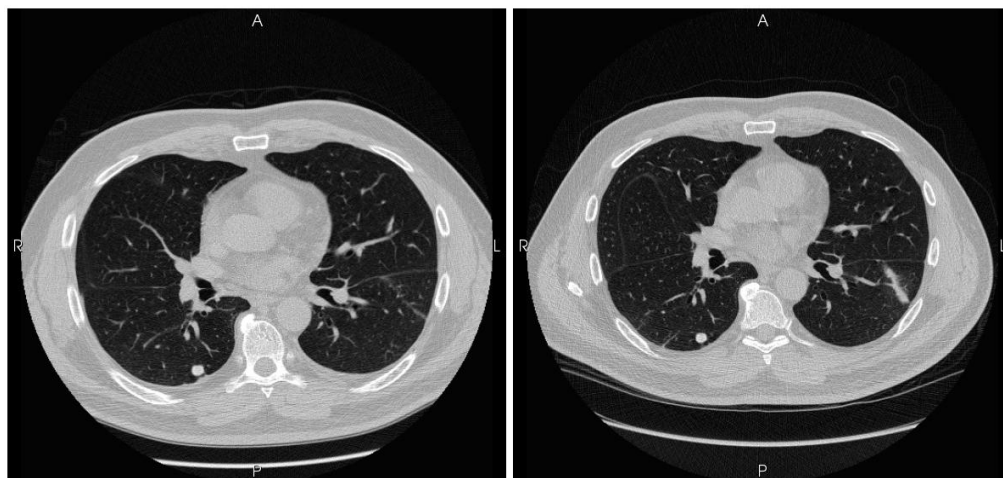


Figure 9. Two scans of a lesion in the VOLCANO Dataset ([View 7](#)).

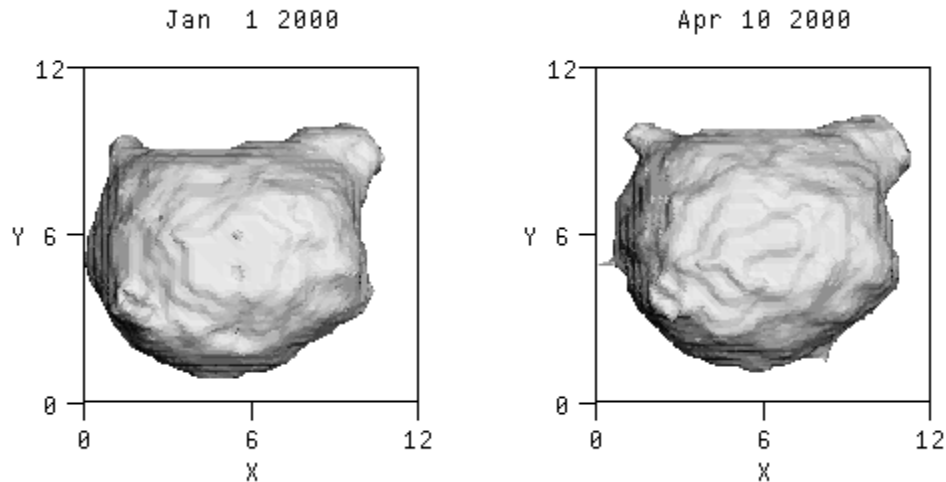


Figure 10: An example of computer assisted segmentation for the lesions shown in Figure 9

The initial findings of this study showed there was no statistical difference between the automated methods on scans of the same slice thickness, but there was a statistical difference in the methods when the scan slice thickness is changed (for subgroup B above). The behavior of the methods for nodules with a small real change in size was similar to that for the zero-change data. The last point has implications for the validity of using zero-size change datasets for evaluating nodule measurement performance. There was an interesting concordance between the different automated methods for a measured change in size for some cases in the zero-change dataset. A follow on to this study is VOLCAMAN'10,²⁴ which enlists a number of physicians using simple manual image marking tools to measure the change in size of a subset of the cases used in VOLCANO'09. In this way the variation of experts for the same task will be established and comparisons with computer methods can be made

Discussion

These examples are only a small portion of what could be done to advance the field. Whether considered from the vantage point of providing an objective basis on which to evaluate the relative performance of different candidate methods, or to allow individual groups access to larger data sets than they would otherwise be able to afford individually, or as a primary driver in the effort to harness the strength of current and new technology towards clinically relevant problems, there is a recurrent theme of the importance of public data resources. Moreover, the ability to evaluate the same data in different ways is arguably not only helpful, but in fact necessary, to establish an objective basis for performance assessment.

This paper identifies several early programs to collect and utilize data either directly in the public domain or easily accessible to teams that demonstrate their need for it to consortia or other groups that recognize a role in collecting and curating such data. Likewise, it is published using the nascent method referred to by this journal as "interactive science publishing," which further encourages a means by which not only the results but also the data used in deriving those results is available for public peer review. We support the editors position that such capabilities will not only move the state of the art in scientific publication forward, but the science itself will benefit as more access is granted to independent reviewers. Such capability is concordant with the goals of our group and we are pleased to be able to exercise it for our present purposes.

403 Other working material of the team is maintained on a Wiki page that enables the group activity.²⁵

404 **Acknowledgements**

405 Many members of our Technical Committee have made substantial contributions to the work described
406 in this report. We would like to specifically acknowledge the following people in addition to the authors
407 (alphabetically): Robert Ford (RadPharm), David Gustafson (Intio), Philip F. Judy (NLST), Matthias Thorn
408 (Siemens Healthcare), James Mulshine (Rush), Daniel Nicolson (Definiens), Kevin O'Donnell (Toshiba),
409 and Daniel C. Sullivan (Duke University and Scientific Advisor to RSNA). Linda Bresolin, Joseph Koudelik,
410 and Fiona Miller of the RSNA organized each meeting, kept our records, and built the contents of our
411 wiki.

412 We would further like to acknowledge the support of the NCI Cancer Imaging Program, its RIDER
413 Project²⁶ and the CDRH/NIBIB Laboratory for the Assessment of Medical Imaging Systems (LAMIS).²⁷

414 **References**

- 1 Buckler A. MEDICAL IMAGING CONTINUUM: Path Forward for Advancing the Uses of Medical Imaging in the Development of New Biopharmaceutical Products
- 2 Buckler A., et al., The Use of Volumetric CT as an Imaging Biomarker in Lung Cancer, Academic Radiology, Vol 17, No 1, January 2010, 100-106
- 3 Buckler A., et al., Volumetric CT in Lung Cancer: An Example for the Qualification of Imaging as a Biomarker, Academic Radiology, Vol 17, No 1, January 2010, 107-115
- 4 Clarke LP, Sriram RD, Schilling LB Imaging as a Biomarker: Standards for Change Measurements in Therapy: Workshop summary. Acad Radiol. 2008 Apr; 15(4):501-30. PMID: 18802422 [PubMed - in process]
- 5 [McLennan G, Clarke L P, and Hohl R.](#) Imaging as a biomarker for therapy response: cancer as a prototype for the creation of research resources. Clin Pharmacol Ther. 2008 Oct; 84(4):433-6. PMID: 18802422 [PubMed - in process]
- 6 Armato S 3rd, Meyer C, McNitt-Gray M, McLennan G, Reeves A, Croft B, Clarke L. The Reference Image Database to Evaluate Response to Therapy in Lung Cancer (RIDER) Project: A Resource for the Development of Change-Analysis Software. Clin Pharmacol Ther. 2008 Oct; 84(4):448-56. Epub 2008 Aug 27. PMID: 18754000 [PubMed - in process]
- 7 Petrick N, Brown DG, Suleiman O, and Myers KJ, Imaging as a Tumor Biomarker in Oncology Drug Trials for Lung Cancer: The FDA Perspective. Clin Pharmacol Ther. 2008, Oct 84 (4): 523-5.
- 8 Goodsitt MM, Chan H-P, Way TW, Larson SC, Christodoulou EG, Kim J. Accuracy of the CT numbers of simulated lung nodules imaged with multi-detector CT scanners. Medical Physics 2006; 33:3006-3017.
- 9 Gavrielides MA, Kinnard LM, Myers KJ, et al. A resource for the development of methodologies for lung nodule size estimation: Database of thoracic CT scans of an anthropomorphic phantom. Optics Express 2010; (Submitted to this special issue).

-
- ¹⁰ Nioutsikou E, Richard N Symonds-Taylor J, Bedford JL, Webb S. Quantifying the effect of respiratory motion on lung tumour dosimetry with the aid of a breathing phantom with deforming lungs. *Physics in Medicine & Biology* 2006; 51:3359-3374, Electronic link: <http://iopscience.iop.org/0031-9155/51/14/005/?jredirect=.iopscience>, accessed: 28-Jan, 2010..
- ¹¹ Way TW, Chan H-P, Goodsitt MM, et al. Effect of CT scanning parameters on volumetric measurements of pulmonary nodules by 3D active contour segmentation: a phantom study. *Physics in Medicine & Biology* 2008; 53:1295-1312.
- ¹² Gavrielides MA, Kinnard LM, Myers KJ, et al. Noncalcified lung nodules: volumetric assessment with thoracic CT. *Radiology* 2009; 251:26-37.
- ¹³ Gavrielides MA, Zeng R, Kinnard LM, Myers KJ, Petrick N. A template-based approach for the analysis of lung nodules in a volumetric CT phantom study. In: *Medical Imaging 2009: Computer-Aided Diagnosis*. 1 ed. Lake Buena Vista, FL, USA: SPIE, 2009; 726009-726011
- ¹⁴ McNitt-Gray MF, Bidaut LM, Armato SG, Meyer CR, Gavrielides MA, Fenimore C, McLennan G, Petrick N, Zhao B, Reeves AP, Beichel R, Kim HJ, Kinnard L., Computed tomography assessment of response to therapy: tumor volume change measurement, truth data, and error. *Transl Oncol*. 2009 Dec; 2(4):216-22.PMID: 19956381.
- ¹⁵ Meyer CR, Armato SG, Fenimore CP, McLennan G, Bidaut LM, Barboriak DP, Gavrielides MA, Jackson EF, McNitt-Gray MF, Kinahan PE, Petrick N, Zhao B. Quantitative imaging to assess tumor response to therapy: common themes of measurement, truth data, and error sources. *Transl Oncol*. 2009 Dec;2(4):198-210.PMID: 19956379.
- ¹⁶ Zhao B, James LP, Moskowitz CS, Guo P, Ginsberg MS, Lefkowitz RA, Qin Y, Riely GJ, Kris MG, Schwartz L. Evaluating variability in tumor measurements from same-day repeat CT scans of patients with non-small cell lung cancer. *Radiology*. 2009 Jul;252(1):263-72.
- ¹⁷ Anthony P. Reeves, Alberto M. Biancardi, , David F. Yankelevitz, Sergei Fotin, Brad M. Keller, Artit Jiraptnakul, and Jaesung Lee. A public image database to support research in computer aided diagnosis. In *31st Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pages 3715-3718, Sept. 2009.
- ¹⁸ <http://www.itl.nist.gov/iad/894.05/biochange2008/Biochange2008-webpage.htm>, last visited 27 January 2010.
- ¹⁹ A. P. Reeves, A. C. Jirapatnakul, A. M. Biancardi, T. V. Apanasovich, C. Schaefer, J. J. Bowden, M. Kietzmann, R. Korn, M. Dillmann, Q. Li, J. Wang, J. H. Moltz, J. Kuhnigk, T. Hayashi, X. Zhou, H. Fujita, T. Duindam, B. van Ginneken, R. Avila, J. P. Ko, K. Melamud, H. Rusinek, R. Wiemker, G. Soza, C. Tietjen, M. Thorn, M. F. McNitt-Gray, Y. Valenciaga, M. Khatonabadi, Y. Kawata, and N. Niki. "The VOLCANO'09 challenge: Preliminary results," In *Second International Workshop of Pulmonary Image Analysis*, pp. 353-364, Sept. 2009.
- ²⁰ A. P. Reeves, A. M. Biancardi, D. Yankelevitz, S. Fotin, B. M. Keller, A. Jirapatnakul, J. Lee. "A Public Image Database to Support Research in Computer Aided Diagnosis," In *31st Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pp. 3715-3718, Sept. 2009.
- ²¹ <http://www.via.cornell.edu/databases/crpf.html> Public Database to Address Drug Response.

-
- ²² <http://preventcancer.org/> The Prevent Cancer Foundation.
- ²³ Reeves, A. P., Biancardi, A. M., Apanasovich, T. V. et al.: The Lung Image Database Consortium (LIDC): A Comparison of Different Size Metrics for Pulmonary Nodule Measurements. Academic Radiology 14, 1475--1485 (2007).
- ²⁴ <http://www.via.cornell.edu/volcaman/> Draft version of the VOLCAMAN study.
- ²⁵ http://qibawiki.rsna.org/index.php?title=Volumetric_CT.
- ²⁶ <http://imaging.cancer.gov/reportsandpublications/ReportsandPresentations/LungImaging/print>.
- ²⁷ <http://www.nibib.nih.gov/Research/Intramural/LAMIS>.