# TABLES

**Table 1.** Demographics, cognitive status and vascular risk factors of participant groups categorized by diagnosis

| | Total (n=45) | Alzheimer disease (n=15) | Posterior cortical atrophy (n=5) | Dementia with Lewy body (n=5) | Frontotemporal Dementia (n=5) | Vascular cognitive impairment (n=5) | Non-neurodegenerative (n=10) |
|---|---|---|---|---|---|---|---|
| **Age, years (mean±sd)*** | 69.44±7.9 | 72.2±8.7 | 69.4±8.6 | 71.83±5.4 | 66±7.1 | 69.8±8 | 65.7±7.1 |
| **Sex** | | | | | | | |
| **Men (n(%))** | 22(48.9) | 9(60) | 1(10) | 2(40) | 4(40) | 3(60) | 3(30) |
| **Women (n(%))** | 23(51.1) | 6(40) | 4(40) | 3(60) | 1(10) | 2(40) | 7(70) |
| **Race** | | | | | | | |
| **Caucasian (n(%))** | 41(91.1) | 14(93.3) | 5(100) | 4(80) | 5(100) | 4(80) | 9(90) |
| **African American (n(%))** | 4(8.9) | 1(6.7) | - | 1(20) | - | 1(20) | 1(10) |
| **Education, years** | | | | | | | |
| **Grade school (n(%))** | 2(4.3) | 1(6.7) | 1(20) | 1(20) | - | 1(20) | - |
| **High school (n(%))** | 16(34.8) | 4(26.7) | - | 1(20) | 3(60) | - | 2(20) |
| **College (n(%))** | 14(30.4) | 3(20) | 2(40) | 3(60) | 2(40) | 2(40) | 6(60) |
| **Post-graduate (n(%))** | 9(19.6) | 6(40) | 1(20) | - | - | 2(40) | 2(20) |
| **Other (n(%))** | 4(8.75) | 1(6.7) | 1(20) | - | - | - | - |
| **Duration of symptoms, years (median(Q1-Q3)*** | 3(2-4) | 3(2-4) | 3(1.5-7.5) | 3(2-6) | 2.7(1.7-5.6) | 5(2-6.5) | 2.5(1-4) |
| **BMI, kg/m2 (mean±sd)** | 26.5±5.7 | 24.7±3.7 | 21.7±4.1 | 25.4±7 | 32.2.±6 | 30.5±5.1 | 27.3±5.8 |
| **CDR global score** | | | | | | | |
| **0 (n(%))** | 5(11.1) | - | - | - | - | 1(20) | 4(40) |
| **0.5 (n(%))** | 25(55.6) | 8(53.3) | 4(40) | 2(40) | 4(80) | 3(60) | 4(40) |
| **1 (n(%))** | 11(24.4) | 6(40) | - | 2(40) | 1(20) | 1(20) | 1(10) |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **2 (n(%))** | 4(8.9) | 1(6.7) | 1(10) | 1(20) | - | - | 1(10) |
| **CDR-SB** (median(Q1-Q3)*) | 2.5(1.5-5.5) | 4.5(2.5-5.5) | 1.5(1.2-8) | 5(2-9) | 3(1-3.5) | 2.5(1-4.2) | 1(0-2.5) |
| **MMSE score** (median(Q1-Q3)*) | 25(20-28) | 21(17-26) | 26(14-26.5) | 19(15-25) | 28(26-28.5) | 28(25.5-29) | 28(23-30) |
| **GDS score** | | | | | | | |
| **0-4 (n(%))** | 35 | 13 | 4 | 5 | 3 | 3 | 7 |
| **5 and above (n(%))** | 10 | 2 | 1 | 0 | 2 | 2 | 3 |

_Abbreviations:_ *BMI: body-mass index;  CDR: Clinical Dementia Rating Scale; MMSE: Mini Mental-State Examination; CDR-SB: Clinical Dementia Rating Scale sum of boxes; PACC:*

*Preclinical Alzheimer Cognitive Composite; GDS: 15-item Geriatric Depression Scale*

*\* Data is presented as mean and standard deviation (mean±sd) when the variable followed a normal distribution and as median plus the first and third quartiles (Q1-Q3) where variables had a none normal*

*distribution. Normality was determined using the Kolmogrov-Smirnov Goodness-of-Fit test*

**Table 2.** Comparison of findings in the final report between ILP and AIRC tool

| | | | ILP based on FreeSurfer 7.1.1 | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Frontal lobe thickness | | Occipital lobe thickness | | Parietal lobe thickness | | FTD thickness* | | Hippocampal volume | | Lateral ventricle volume | | Ventricle/Cerebrum ratio* | |
| | | | >-2SD | <-2SD | >-2SD | <-2SD | >-2SD | <-2SD | >-2SD | <-2SD | >-2SD | <-2SD | >+2SD | >+2SD | <+2SD | >+2SD |
| AIRC Tool | Frontal lobe volume* | >10% | 28 | 3 | | | | | | | | | | | | |
| | | <10% | 2 | 12 | | | | | | | | | | | | |
| | Occipital lobe volume** | >10% | | | 36 | 2 | | | | | | | | | | |
| | | <10% | | | 3 | 4 | | | | | | | | | | |
| | Parietal lobe volume** | >10% | | | | | 28 | 1 | | | | | | | | |
| | | <10% | | | | | 1 | 15 | | | | | | | | |
| | Frontal and temporal lobe volumes§ | >10% | | | | | | | 31 | 5 | | | | | | |
| | | <10% | | | | | | | 0 | 9 | | | | | | |
| | Hippocampal volume** | >10% | | | | | | | | | 34 | 2 | | | | |
| | | <10% | | | | | | | | | 2 | 7 | | | | |
| | Lateral ventricle volume§§ | <90% | | | | | | | | | | | 29 | 0 | | |
| | | >90% | | | | | | | | | | | 4 | 12 | | |
| | | <90% | | | | | | | | | | | | | 29 | 0 |

| | | | | | |
|---|---|---|---|---|---|
| | Lateral ventricle volume$^{§§}$ | **>90%** | | **3** | 13 |

*Abbreviations*: ILP: individual longitudinal participant, FTD: frontotemporal dementia

* FTD thickness: average thickness and frontal and temporal region volumes implicated in FTD; Ventricle/Cerebrum ratio: see Supplementary Table 1

** demonstrates number of participants with/without frontal, occipital, parietal lobe or hippocampal normalized volumes atrophy (below and above 10[th] percentile compared to normative data) in *at least one hemisphere* based on the final report generated by the AIRC brain morphometry tool

§ demonstrates number of participants with/without *both* frontal and temporal lobe atrophy (below and above 10[th] percentile compared to normative data) in *at least one hemisphere* based on the final report generated by the AIRC brain morphometry tool.

§§ demonstrates number of participants with/without lateral ventricle enlargement (above and below 90th percentile) in *at least one hemisphere* based on the final report generated by the AIRC brain morphometry tool.

# SECTION 1. ATLAS/REGION MATCHED AND SUMMARY STATISTICS IN THE FREESURFER/ILP PIPELINE

**Supplementary Table 1.** Summary metrics generated by the Individual Longitudinal Participant (ILP) tool and their constituent FreeSurfer outputs

| Frontal Lobe Thickness§ | FTD Thickness§ | Ratio of Lateral Ventricle to Cerebral volume | Occipital Lobe Thickness§ | Parietal Lobe Thickness§ |
|---|---|---|---|---|
| | | Cerebral volume is the supratentorial volume minus following structures: | | |
| L&R superior frontal thickness | Lateral orbitofrontal thickness | Third ventricle | L&R lateral occipital thickness | L&R postcentral thickness |
| L&R rostral middle frontal thickness | Pars opercularis thickness | Fourth ventricle | L&R cuneus thickness | L&R supramarginal thickness |
| L&R caudal middle frontal thickness | Rostral middle frontal thickness | Fifth ventricle | L&R pericalcarine thickness | L&R superior parietal thickness |
| L&R precentral thickness | Pars triangularis thickness | L&R inferior lateral ventricle | L&R lingual thickness | L&R inferior parietal thickness |
| L&R pars opercularis thickness | Insula thickness | L&R lateral ventricle | | |
| L&R pars triangularis thickness | Inferior temporal thickness | L&R choroid plexus | | |
| L&R pars orbitalis thickness | Middle temporal thickness | CSF | | |
| L&R lateral orbitofrontal thickness | | | | |
| L&R medial orbitofrontal thickness | | | | |

§ calculated as the average thickness of the constituent cortical regions
L: left hemisphere; R; right hemisphere

**Supplementary Table 2.** Approximate mapping of Desikan-Killiany FreeSurfer-based regions of interest to different brain lobes from the AIRC output

| Frontal lobe [§] | Occipital lobe[§] | Temporal lobe[§] | Parietal lobe[§] | Cingulate[§] | Cerebellum |
|---|---|---|---|---|---|
| Caudal middle frontal | Cuneus | Banks of superior temporal sulcus | Inferior parietal | Caudal anterior cingulate | Cerebellum white matter |
| Frontal pole | Lateral occipital | Entorhinal | Postcentral | Isthmus cingulate | Cerebellum cortex |
| Lateral orbitofrontal | Lingual | Fusiform | Precuneus | Posterior cingulate | |
| Medial orbitofrontal | Pericalcarine | Inferior temporal | Superior parietal | Rostral anterior cingulate | |
| paracentral | | Middle temporal | Supramarginal | | |
| Pars opercularis | | Parahippocampal | | | |
| Pars orbitalis | | Superior temporal | | | |
| Pars triangularis | | Temporal pole | | | |
| Precentral | | Transverse temporal | | | |
| Rostral middle frontal | | | | | |
| Superior frontal | | | | | |

§ The sum of volumes of the constituent FreeSurfer regions was calculated and used in subsequent comparisons with the lobar volumes from AIRC Brain MR tool

# SECTION 2. OPERATIONAL TERMS AND PROTOCOLS

**A: Correlation, Agreement and Consistency**:

1. ***Pearson's Correlation Coefficient (PCC)***: PCC is a measure of linear correlation between two sets of data, measured as the ration between the covariance of the two variable over their pooled standard deviations. A Pearson correlation can be a valid estimator of interrater reliability, but only when meaningful pairings exist between two and only two raters, making PCC less sensitive to inter-rater bias. Other limitations of PCC are that it is unable to reliably detect nonlinear relationships and is very sensitive to outliers.

2. ***Intraclass Correlation Coefficient (ICC)***: ICC is a descriptive statistic of quantitative measurements that are based on units (i.e. scans) categorized into groups (e.g volumetric software). ICC describes how strongly units in each group resemble each other. ICC is in that sense similar to PCC as both show how similar two sets of measurements are to each other. Remember however, that PCC is unable to reliably measure inter-rater bias, as each variable is centered and scaled by its own mean and standard deviation. In ICC however the data are centered and scaled using a pooled mean and standard deviation. There are two ways by which ICC can be modeled: 1) *Absolute Agreement*, and 2) *Consistency*. In the context of different measurements by two independent raters (here volumetric software), absolute agreement reflects whether different raters assign the same measurement to the same subject. Conversely, consistency is sensitive

whether measurements by different raters to the same group of subjects are correlated in an additive manner. As a result:

i. ***ICC-consistency***: This measurement is blind to systematic differences (or errors) between raters and only takes into account the random residual error. ICC consistency is therefore not sensitive enough to detect whether one software systematically over or under estimates the volume of a given brain area, compared to the other or gold standard software

ii. ***ICC-agreement***: accounts for both systematic errors of both raters and random residual errors. As a result ICC agreement is often lower than ICC consistency.

**B: Compatibility**:

We labeled the radiologist impression compared to the established clinical diagnosis as either *compatible* or *non-compatible*.

Each patient received a clinical diagnoses by a skilled neurologist, at the end of the clinical visit and prior to any neuroimaging. These were one of the following categories:

A) ***Neurodegenerative***: Alzheimer disease (AD), posterior cortical atrophy (PCA), dementia with Lewy bodies (DLB), frontotemporal dementia (FTD) and vascular cognitive impairment (VCI). Scans picked within any of the above categories were all showing evidence of symmetric, lobar neurodegeneration (whether grossly in structural MRI, or more likely within volumetric quantifications).

B) ***Non-Neurodegenerative***: Conditions such as subjective cognitive impairment in the absence of clinical dementia, mood disorders, polypharmacy and sleep disorders causing the cognitive symptoms. These scans had no structural abnormality.

## C: Radiologist Rater Assessment for Compatibility

Radiologists were only informed of the participant's age and sex and that he/she was being assessed for a chief complaint of cognitive impairment. We compared the radiologists answer to the following questions to the known disease entity or category from the above:

Q1. Are there any abnormal findings in the study 'suggestive of a cause for dementia? (Yes/No)

Q2. Are the findings symmetric or asymmetric? (Yes/No)

Q3. Is there evidence of lobar atrophy? (Yes/No)

Q4. Do findings point to a specific neurodegenerative entity? (Yes/No)

Q5. If yes, what is your clinical impression (multiple choices: AD, PCA, VCI, FTD, or DLB)

If the patient clinical diagnosis was a ***Neurodegenerative*** condition, the radiologist's impression was considered ***Compatible*** if the radiologist:

i.      Answered Yes to Q1 and at least two of the Q2-Q4, <u>or;</u>

ii.     Answered Yes to Q1 and picked the correct diagnosis among choices in Q5

If the patient clinical diagnosis was a ***Non-neurodegenerative*** condition, the radiologist's impression was considered ***Compatible*** if the radiologist:

i.       Answered No to Q1 - this would automatically end the evaluation and rest of question would not be demonstrated, <u>or;</u>

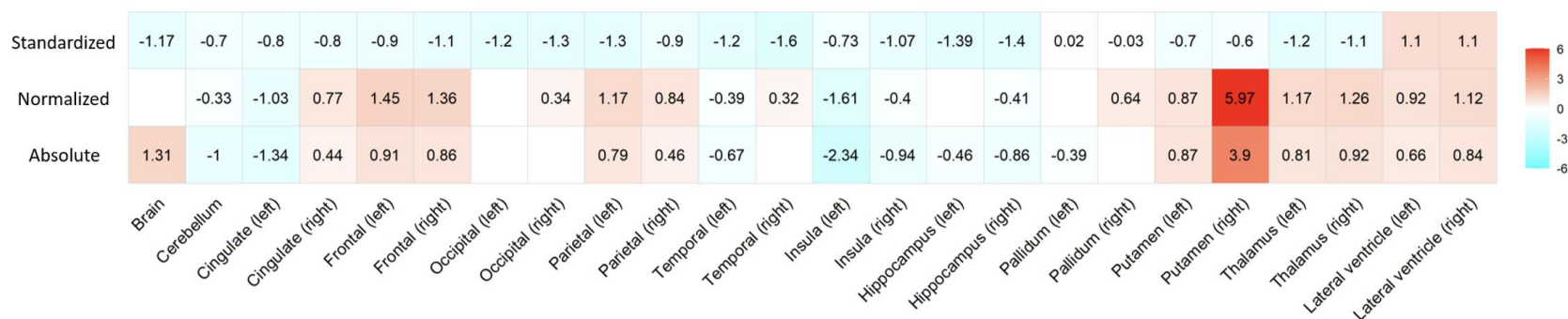ii.      Answered Yes to Q1, No to at least one of the Q2 or Q3, and No to both Q4 and Q5

# SECTION 3. PAIRED STATISTICS

In order to compare the volumetric values between FS and AIRC tools and extract the mean differences in volumes and effect sizes, we further performed either a paired-sample t-test using the "t.test" function (if normality assumption was met) or a paired-samples Wilcoxon signed-rank test using the "wilcox-test" function (if normality assumption was not met). An effect size of 0.2 or below was interpreted as negligible difference in the measurements while the range of 0.2-0.5, 0.5-0.8 and >0.8 indicated small, medium and large effect sizes, respectively [1].

**A) Absolute volumes**

When compared through paired statistics, absolute volumes measured by the AIRC tool and FS were significantly different, except in the bilateral occipital lobes, the right temporal lobe and the right pallidum. The respective effect sizes were medium to large for these significant regions (Supplementary Figure 1, and Supplementary Table 3).

**Supplementary Figure 1.** Comparing effect size of paired t-test statistics between volumetric measurements produced by the AIRC versus FS/ILP tools



*Panels demonstrate correlation coefficients for raw volumes (bottom), volumes normalized to TIV (middle) and standardized (z-score) volumes (top). Negative values indicate statistically significant lower volumes calculated by the AIRC compared to the FS/ILP and positive values indicate statistically significant higher volumes in AIRC relative to the FS/ILP. AIRC: AI-Rad Companion brain MR tool; FS/ILP: FreeSurfer/Individual Longitudinal Participant Pipeline; TIV: total intracranial volume*

**Supplementary Table 3.** Paired T-test comparing absolute regional volumes between FS 7.1.1 and AIRC Brain MR tool

| | P-value* | Mean Difference**(95%CI) (ml$^3$) | Effect Size$^\S$ |
|---|---|---|---|
| Total Intracranial volume | <0.001 | -186.8((-290.2)–(-83.5)) | -0.543 |
| Brain | <0.001 | -64.8 ((-79.7)–(-49.9)) | 1.309 |
| Cerebellum | <0.001 | -10.5 ((-13.6)–(-7.3)) | -1.000 |
| Cingulate (left) | <0.001 | -1.9((-2.3)–(-1.5)) | -1.336 |
| Cingulate (right) | 0.006 | 0.5(0.18-0.9) | 0.441 |
| Frontal (left) | <0.001 | 8.3(5.6-11.1) | 0.914 |
| Frontal (right) | <0.001 | 7.9(5.1-10.7) | 0.859 |
| Occipital (left) | 0.17 | -0.6((-1.4)–(0.2)) | -0.221 |
| Occipital (right) | 0.97 | -0.15((-0.9)–(0.9)) | -0.005 |
| Parietal (left) | <0.001 | 5.7(3.5-7.9) | 0.792 |
| Parietal (right) | 0.005 | 3.1(1.0-5.1) | 0.455 |
| Temporal (left) | <0.001 | -4.2((-6.1)–(-2.3)) | -0.667 |
| Temporal (right) | 0.93 | -0.1((-2.1)–(1.8)) | -0.019 |
| Insula (left) | <0.001 | -1.5((-1.7)–(-1.3)) | -2.335 |
| Insula (right) | <0.001 | -0.5((-0.7)–(-0.4)) | -0.935 |
| Hippocampus (left) | 0.004 | -0.2((-0.3)–(-0.082)) | -0.461 |

| | | | |
|---|---|---|---|
| Hippocampus (right) | **<0.001** | -0.4((-0.5)–(-0.2)) | -0.865 |
| Pallidum (left) | **0.016** | -0.09((-0.1)–(-0.02)) | -0.388 |
| Pallidum (right) | 0.55 | 0.02((-0.04)–(0.08)) | 0.100 |
| Putamen (left) | **<0.001** | 2.6(2.5-2.8) | 0.87 |
| Putamen (right) | **<0.001** | 2.1(1.9-2.2) | 3.902 |
| Thalamus (left) | **<0.001** | 0.7(0.4-0.9) | 0.812 |
| Thalamus (right) | **<0.001** | 0.8(0.5-1.1) | 0.915 |
| Lateral ventricle (left) | **<0.001** | 2.3(1.2-3.3) | 0.660 |
| Lateral ventricle (right) | **<0.001** | 2.6(1.7-3.6) | 0.838 |

*p-values from either a paired t-test (all regions except for left putamen) or Wilcoxon signed rank statistics (left putamen) are corrected for multiple comparison using the Benjamini-Hochberg FDR correction method.*

**Mean Differences and their 95%confidence intervals (95%CI) represent the structures volume in FreeSurfer minus AIRC Brain MR tool calculated through a paired t-test (all regions except for left putamen) or Wilcoxon signed rank statistics (left putamen) depending on the normality assumption.*

§ *Effect size represents either the Cohen's D (all regions except for left putamen) or a non-parametric effect size calculated through the r=z/sqrt(n) formula (left putamen). Effect size estimates from tests that reached statistical significance are color coded based on the magnitude of effect size as below:*

| Negligible<0.2 | Small: 0.2-0.5 | Medium: 0.5-0.8 | Large >0.8 |
|---|---|---|---|

**B) Normalized Volumes**

Normalization of the absolute volumes resulted in a small to moderate effect size when volumes were compared using paired statistics between the two tools. Paired t-test demonstrated a large difference in the normalized volumes of the bilateral frontal lobes and lateral ventricles, while differences in the normalized volumes of the occipital, parietal and temporal lobes and the bilateral hippocampi were small to non-existent between the two tools (Supplementary Figure 1, and Supplementary Table 3).

**Supplementary Table 4.** Paired T-test comparing normalized regional volumes between FS 7.1.1 and AIRC Brain MR tool

| | P-value* | Mean Difference**(95%CI) (%TIV) | Effect Size[§] |
|---|---|---|---|
| Brain | 0.333 | -0.16((-0.61)-(0.9)) | -0.65 |
| Cerebellum | **0.036** | -0.26((-0.51)-(-0.026)) | -0.333 |
| Cingulate (left) | **<0.001** | -0.1((-0.14)-(-0.07)) | -1.034 |
| Cingulate (right) | **<0.001** | 0.07(0.04-0.09) | 0.767 |
| Frontal (left) | **<0.001** | 0.8(0.65-1) | 1.453 |
| Frontal (right) | **<0.001** | 0.8(0.6-1) | 1.362 |
| Occipital (left) | 0.37 | 0.025((-0.027)-(0.07)) | 0.144 |
| Occipital (right) | **0.03** | 0.068(0.01-0.13) | 0.342 |
| Parietal (left) | **<0.001** | 0.54(0.4-0.68) | 1.171 |
| Parietal (right) | **<0.001** | 0.36(0.23-0.49) | 0.838 |
| Temporal (left) | **0.01** | -0.15((-0.27)-(-0.03)) | -0.394 |
| Temporal (right) | **0.04** | 0.124(0.006-0.24) | 0.315 |
| Insula (left) | **<0.001** | -0.08((-0.1)-(-0.06)) | -1.613 |
| Insula (right) | **0.01** | -0.018((-0.032)-(-0.0045)) | -0.399 |
| Hippocampus (left) | 0.52 | -0.004((-0.01)-(0.007)) | -0.102 |
| Hippocampus (right) | **0.01** | -0.01((-0.025)-(-0.003)) | -0.409 |

| | | | |
|---|---|---|---|
| Pallidum (left) | 0.735 | 0.0008((-0.004)-(0.057)) | 0.051 |
| Pallidum (right) | **<0.001** | 0.009(0.004-0.013) | 0.640 |
| Putamen (left) | **<0.001** | 0.2(0.19-0.21) | 0.87 |
| Putamen (right) | **<0.001** | 0.16(0.15-0.17) | 5.970 |
| Thalamus (left) | **<0.001** | 0.07(0.05-0.09) | 1.171 |
| Thalamus (right) | **<0.001** | 0.08(0.06-0.1) | 1.258 |
| Lateral ventricle (left) | **<0.001** | 0.25(0.16-0.3) | 0.921 |
| Lateral ventricle (right) | **<0.001** | 0.26(0.19-0.33) | 1.119 |

_Abbreviations:_ _TIV: total intracranial volume_

_*p-values from either a paired t-test (all regions except for left putamen) or Wilcoxan signed rank statistics (left putamen) are corrected for multiple comparison using the Benjamini-Hochberg FDR correction method._

_**Mean Differences and their 95% confidence intervals (95%CI) represent the structures volume in FreeSurfer minus AIRC Brain MR tool calculated through a paired t-test (all regions except for left putamen) or Wilcoxan signed rank statistics (left putamen) depending on the normality assumption._

_§ Effect size represents either the Cohen's D (all regions except for left putamen) or a non-parametric effect size calculated through the r=z/sqrt(n) formula (left putamen). Effect size estimates from tests that reached statistical significance are color coded based on the magnitude of effect size as below:_

| | | | |
|---|---|---|---|
| Negligible<0.2 | Small: 0.2-0.5 | Medium: 0.5-0.8 | Large >0.8 |

**C) Standardized Volumes**

When assessed through paired t-test statistics, the standardized volumes were largely different in all cortical regions, bilateral hippocampi and the lateral ventricles, when the AIRC and FS/ILP outputs were compared. Z-scores calculated by the FS/ILP pipeline were between 0.8 to 1.8 scores lower in the main four cortical regions and the hippocampi compared to the AIRC tool, corresponding to a large effect size in these regions (Figure 3 and supplementary Table 5). Given the non-significant to small difference in the normalized volumes, this might owe to differences in the composition of the normative datasets for each cohort.

**Supplementary Table 5.** Paired T-test comparing regional z-scores between FS 7.1.1 and AIRC Brain MR tool

| | P-value* | Mean Difference**(95%CI) (%TIV) | Effect Size[§] |
|---|---|---|---|
| Brain | **<0.001** | -1.31((-1.65)-(-0.97)) | -1.17 |
| Cerebellum | **<0.001** | -0.7((-1)-(-0.4)) | -0.7 |
| Cingulate (left) | **<0.001** | -1.17((-1.62)-(-0.7)) | -0.8 |
| Cingulate (right) | **<0.001** | -1.01((-1.4)-(-0.62)) | -0.8 |
| Frontal (left) | **<0.001** | -1.2((-1.57)-(-0.84)) | -0.9 |
| Frontal (right) | **<0.001** | -1.32((-1.68)-(-0.96)) | -1.1 |
| Occipital (left) | **<0.001** | -1.23((-1.52)-(0.93)) | -1.2 |
| Occipital (right) | **<0.001** | -1.48((-1.8)-(1.16)) | -1.3 |
| Parietal (left) | **<0.001** | -1.4((-1.83)-(-0.97)) | -1.3 |
| Parietal (right) | **<0.001** | -1.4((-1.74)-(-1.07)) | -0.9 |
| Temporal (left) | **<0.001** | -1.25((-1.59)-(-0.9)) | -1.2 |
| Temporal (right) | **<0.001** | -1.64((-2.02)-(-1.26)) | -1.6 |
| Insula (left) | **<0.001** | -0.5(-0.83)-(-0.15)) | -0.73 |
| Insula (right) | **<0.001** | -0.87((-1.2)-(-0.51)) | -1.07 |
| Hippocampus (left) | **<0.001** | -1.26((-1.63)-(-0.89)) | -1.39 |
| Hippocampus (right) | **<0.001** | -1.21((-1.56)-(-0.86)) | -1.4 |

| | | | |
|---|---|---|---|
| Pallidum (left) | 0.07 | 0.047((-0.54)-(0.64)) | 0.024 |
| Pallidum (right) | **0.04** | -0.35((0.68)-(-0.019)) | -0.03 |
| Putamen (left) | **<0.001** | -0.47((-0.69)-(-0.26)) | -0.7 |
| Putamen (right) | **<0.001** | -0.5((-0.7)-(-0.29)) | -0.6 |
| Thalamus (left) | **<0.001** | -1.37((-1.73)-(-1.02)) | -1.2 |
| Thalamus (right) | **<0.001** | -1.4((-1.79)-(-1.02)) | -1.1 |
| Lateral ventricle (left) | **<0.001** | 0.7(0.49-0.91) | 1.1 |
| Lateral ventricle (right) | **<0.001** | 0.82(0.58-1.05) | 1.1 |

*Abbreviations: TIV: total intracranial volume*

*\*p-values from a paired t-test that are corrected for multiple comparison using the Benjamini-Hochberg FDR correction method.*

*\*\*Mean Differences and their 95% confidence intervals (95%CI) represent the structures volume in FreeSurfer minus AIRC Brain MR tool calculated through a paired t-test*

*§ Effect size represents either the Cohen's D. Effect size estimates from tests that reached statistical significance are color coded based on the magnitude of effect size as below:*

| Negligible<0.2 | Small: 0.2-0.5 | Medium: 0.5-0.8 | Large >0.8 |
|---|---|---|---|

# SECTION 4. MISCELLANEOUS

**Supplementary Table 6.** Percentage of correct responses to each question by radiologist categorized by the tool among participants with neurodegenerative diagnosis

| Method | Presence of abnormality? (%correct) | Presence of symmetric & lobar atrophy? (%correct) | Which neurodegenerative entity? (%correct) |
|---|---|---|---|
| MPRAGE_only | 54.3% | 48.4% | 78.26% |
| MPRAGE+ILP | 74.3% | 48.5% | 68% |
| MPRAGE+AIRC | 71.4% | 54.2% | 81.25% |
| *Abbreviations:* ILP: individual longitudinal participant tool; AIRC: AI rad companion software | | | |

**Supplementary Figure 2.** A side-by-side comparison of labeling in FreeSurfer versus AIRC

*Footnote: Panel A and B) White matter is labeled with red in FreeSurfer and light blue in AIRC. Medium blue and dark blue colors in AIRC output designate cortical grey matter and CSF respectively. Note the incorrect labeling of the white matter to grey matter or CSF in the ARIC output. Panel C) hippocampus is labeled with the yellow (FreeSurfer) or green color (AIRC) and is located immediately beneath the inferior horn of the lateral ventricle. Panel D) the globus pallidus (lateral) and putamen (medial) are delineated. In panels D and C, note the relative size of the subcortical structures in FS or AIRC compared to the T1 image. Note the slight difference in location of slices in sagittal view in the FS output, resulting from resampling of the native T1 image to the atlas space during the visualization. This participant had non-neurodegenerative cognitive dysfunction.*