

	Topic	Discussion Notes	TPS Reaction	Action	Thurs. Call Topic	ATR	Future TAC Topic	Future TPS Topic	Quarterly Mtg Topic	Board Meeting	Action for CAL/WIDA PM Team	Action for CAL Psychometrics Team	Action for WIDA
A1	Overview	<p>Questions- Lyle - lacks conceptual framework for guiding - comparability of interpretation Haven't articulated warrants about interpretation Maintain scale scores is mechanism - urge to think of claims and warrants - what kind of evidence required to support Gary - may not want to maintain same scale score - want to make sure scale scores mean the same thing Lyle - interpretation - claim is the scale scores is the same Push beyond looking at scores, articulate warrants about interpretation - Lyle - all studies in service of interpretation Gary - represent same construct Jamal- comparability of PB and computer- accommodation- sources of construct irrelevant variance- cannot use as accommodation Assuming 65% of student will be cb, may be higher. Idea is to move to mostly online as soon as possible. Paper is accommodation - working with APIP to build in accommodation Some states not part of PARC, SBAC, but those are all online states Lyle- in report - provide explanation of adaptations exactly Lyle: Need to address certain questions: o What are the AUA claims for ACCESS for ELLs Online? o What kinds of evidence do we collect to support AUA claims? o Articulate about interpretation</p>	<p>"Back and forth between Gary and Tim: Evidence as an outcome vs. evidence as an input. • Steve's reaction: What evidence do we have that the data is showing the same thing? o Need to re-think old assumptions because the Consortium now has 37 states with a more diverse population. Parameters are not stable. Relying on 350 data points may not reflect this. How broad of a picture do we want to paint? o CAL: The representation of the population from the beginning has grown. Can go back and see how the population has differed from the past. "Anything coming out of the 350 is the initial evaluation. Then look at equating window sample. Shu Jing can send info on the Reading domain Alt ACCESS study. Parameters are robust. • Steve: Consider looking at all domains beyond Reading. Would want this so that the group can feel comfortable saying that the changes in scales is supportable. • Robert: Want a step by step AMAO approach</p>	<p>1. Provide explanation of adaptations for computer based items/tasks in greater detail in report (Lyle's recommendation) - (Alicia and Carsten) send TIDP and Framework doc to Lyle first. Psychometrics team will summarize the adaptation in the ATR 2. See how ACCESS population has differed from the past now that there are 37 states. (no action - new standard setting. How will new states affect the equating? Shu Jing can send a doc and PPT. WIDA will see how we can work with the SEA Relations department) 3. Consider looking at all domains, not just Reading, so TPS group can feel comfortable about saying that changes in scale is supportable. 4. CAL to see how the population has differed from the past now that the Consortium has grown to 37 states. 5. Shu Jing can send info on the Reading domain Alt ACCESS study that shows that parameters are robust to TPS members</p>	<p>*talk about next steps for TPS Quarterly Meeting.</p>			<p>2. Send out one email with follow-up docs and answers to previous meeting before next meeting</p>	<p>*Make sure that TPS suggested agenda topic is something that is helpful to our work.</p>				<p>1. Psychometrics team will summarize the adaptation in the ATR 1. (Alicia and Carsten) send TIDP and Framework doc to Lyle first.</p>
A2.1	Listening - comparability and outliers	<p>Tim- outliers - variation- do we have conceptual model of what outlying data we feel comfortable getting rid of based on mode- and what are effects of removing outliers Last year - TAC recommended we survey students so we did not stack deck - motivated removal Shu Jing - moving to script based - went both ways - look at survey questions- compare student results, - same school, etc. Also, look at different performance, compare to survey- some students perform better, not always worse - not consistent Jamal - to maintain comparability - cb and pb - mode should be the only difference- some other differences beyond just mode (grade clusters changed) - differences which may impact comparability Dorry - specs exist for first grade - look at first grade performance - selected tasks which were appropriate to first grade - no new specs were developed Carol - overall mean performance may change on form with new grade specification Dorry - yes - preparing - looking at difference in performances by grade. Not a norm reference test. Mean performance should fit ability of student - construct has not changed in transition Shu Jing - still measuring same proficiency levels Aki - scale chart - correlation - disattenuated - high approx .8 - with pl, close to 1 thought it was all grade levels combined - quite disperse for a single grade Shu jing - differences in mean - Gary - coefficient of determination not big - comparable charts paper to paper - distributions as part of equating Figure 2 during ft data, one tier- attenuated. figure 3 using operational data - link between paper and computer may not be strong - in order to get more confidence in CB PB, - simulate paper paper correlation and compare to SB- MD change Survey - keep it simple Jamal Are you going to do item level analyses on mode - can do DIF - or Gary - structural model. is there any way to increase the size- Dorry - can use equating sample Gary Analysis would need to be simple - data reduction - are we measuring the same construct</p>		<p>1. Shu Jing provide technical brief of Media Based vs. SB to TAC (no publication date yet). Shu Jing send to WIDA. 2. Look into ways to increase size, perhaps by using equating sample (for Listening). Should we increase the equating sample up from 1000? 3. Look at item invariance 4. address how PB and CB differences impact comparability (like grade cluster change) 5. TAC wants more information on kinds of tasks and differences - specifications in order to get more confidence in CB PB, - simulate paper-paper correlation and compare to SB-MD change. 6. Analysis would need to be simple to determine if we are measuring the same construct across modes - data reduction techniques 7. Also analysis at item level - structural model across mode - see if there is mode effect 8. Carol - having that information can change test specs - look at test spec across modes</p>									
	Listening - Video	<p>Lyle - applied linguistics - comparability of tasks - to what extent is WIDA exploring what CB allows for testing - offering videos - construct is different from face to face conversation- play around with video - push in that direction - TAC wants more information on kinds of tasks WIDA - video - exploring - currently staying as close as possible to what was done. Moving forward introducing</p>		<p>1. Explore video items</p>									
A2.2	Reading - comparability	<p>Lyle - read ahead - also use common item approach - use operationally? Yes Have been using to link pb test - can use baseline of other linking Gary - common item assume variance across mode is consistent Use both - common person and common item - if common item fits, then move on Dorry - specs haven't changed, used older items Gary - folder effects, passage effects Dorry - research done long ago to show they are independent Carol - scrolling? Dorry one item that scrolled, generally no Gary Empirical question - comparability in folders - something different across mode, how students interact with passage Lyle - item types presented in test and item design plan Gary - in peer review w- transition doc - transition is something guidelines will address Jamal - Analyze including and excluding outliers CAL - If there are different results, then you can say that the outliers may have impact on linking. Then what action should we take? This would not be common person equating. Lyle: Compare to common item Gary - outlier is still there - what is decision rule? Lyle- convince stakeholders - Tim - what is belief about outliers? Any data -set 95% cut - will throw out 15 kids in any analysis. Did not finish - strengthened to have some sort of theory for outliers and apply Carol - but is differential performance construct relevant Dorry - conducted cog labs - addressed in different ways- do something systematic. Gary - need to check comparability in the folders. Is there difference in reading on paper vs. computer screen</p>		<p>1. Gary - common item assume variance across mode is consistent Use both - common person and common item - if common item fits, then move on 2. Look at comparability in folders - is there something different across mode, how students interact with passage. 3. peer review guidelines - transition is something guidelines will address- peers will make judgment on comparability of transition based on documentation 4. Do analysis including and excluding outliers - compare- how to choose results? Did outliers have impact on linking? What is decision rule? Have some sort of theory for outliers and apply - motivated - differential performance - systemic</p>									

A2.3	Writing - Scoring +/-	<p>Scoring – Lyle – pluses and minuses take on meaning which may not be similar across proficiency levels – scale should reflect intentions built into ratings Point is not to say it is continuous – make scale reflect interpretations – bigger differences between categories in diff prof levels</p> <p>Score report includes descriptors – scale score and interpretation – generic performance level descriptors from standards – generic across 1-12 in standards – plus what does that mean for diff grade levels – Lyle – wants to see rating scale and interpretive scale – are they similar – meaning should be conveyed in final score</p> <p>Carol – hoping for higher levels of agreement than done in past Dorry – data not yet available</p> <p>Tim - % of scores on access and online access – some indication that raters are selecting some ratings more than other (444) – is that similar in the new scale. That could affect how it is equated. Dorry – focus group with raters this week</p> <p>Lyle – what use is being made of analytic rating vs holistic – Gary – research on teacher report – scores by standard – people do not generally use report – structural model – some relationship, but not used by teachers- Dorry – that is motivation for simplifying Lyle – analytic not being used – then why do it? Dorry – simplify- new raters report this rubric is simpler to use Use questionnaire to inform outlier analysis – effect of keyboarding later Goal – same step parameters for all test forms, but investigate Tim – what rules for scores to be used for linking – what is interrater agreement rules for both – Dorry – look at raw ratings – similar ball park – Tim – necap- all samples double scored – what is done here, Dorry – everything double scored for field test – Gary – pluses and minuses on each task – pilot report to see if it has any meaning Score report based on wide members desire for diagnostic information, but no research done on how to use it Lyle – need to say what pluses and minuses are- could be done for different reasons</p>		<ol style="list-style-type: none"> 1. In report – include interpretive cross walk of score and interpretations 2. In report – include % of scores on access and online access – some indication that raters are selecting some ratings more than other (444) – is that similar in the new scale. That could affect how it is equated. Show how raters are using new rubric 3. Take concern with pluses and minuses to scoring committee – may make raters feel good, but could cause problems with interpretations, Will be presented at next year’s TAC 4. Address concern about application across grade and tiers 5. Try partial credit model on each task 							
A2.3	Writing - Scoring Reliability	<p>Tim – reliability of scores Dorry a.3.2 – linguistic and qualitative analyses – look at differential for paper vs paper. Tim- go back and double score old scores? That is a lot Dorry – double score those that are out of kilter? Contractor only provide one score, but provide statistics on adjudication Carol - Field test scoring – how many raters- everything double scored- best of best- rasch measurement analysis of double scores – raters with different level of severity</p>		<ol style="list-style-type: none"> 1. Do rasch measurement analysis of double scores – raters with different level of severity 							
A2.3	Writing - comparability	<p>Lyle is keyboarding an issue for raters? Yes plans to do those studies Gary – can identify salient differences between ratings – ran structural models, but is form dependent and tier dependent. How do you use it? Lyle – need description of differences – on thin ground Gary how do you know what teachers will do with differences Report for interpretive value- Lyle -Need aua – identify most important claims and provide rebuttal – how to prioritize Tim – control for amount of text? Dorry – there is a maximum online Time Length is not scored – what is language being shown – more should not get higher score- hard to control for – report for comparison with necap- part of scoring rubric – some measures of number of answers before and after In necap - Mode affected scores – kids did better in narrative vs persuasive on keyboard. What are controls for longer prompts on type of writing – Dorry – same specs on ft Gary Composite scale proficiency levels – found was not high enough to look at analyses between content test and access – capture higher end to build scale beyond 5 that is meaningful- to give scores Gary – can find level at which content overtakes language</p> <p>Scale post equating – carol – proposal is to anchor – on both step and task difficulty parameters – can get different results</p>		<ol style="list-style-type: none"> 1. Look at discrimination across spectrum. 2. Jenn – assign her to look at claims for FT 3. conduct analysis on word count of keyboarded vs. handwritten 4. Can we augment FT to make sure we have students at the highest proficiency level? 							
A2.4	Speaking - linking	<p>Correlations – affected by ceiling affect – need to look at correlation of R/L together with speaking Lyle Concern about students being routed to too hard a difficulty level Can elicit information at that low a level? Don’t want to overwhelm student and they don’t know what to do Gary: error at lower level will be greater Carol: it’s now a multi-modal speaking test. You’ve changed the construct of the test. CAL will need to draw out how the delivery of the speaking test is similar and different because the paper speaking test is also in some ways multi-modal with the test administrator and booklet. Treat as a new scale as opposed to restructured scale – equipercentile – but have capped scores to deal with- cannot do this Gary – concern about new test vis a vis AMAO’s – will have big affect if we draw out speaking. Currently some states require no domain lower than 5 – this may change – policy folks will need to make this call – result of getting more consistent and more focused sample of speaking. Examine the impact of classification accuracy Lyle- need to be able to justify differences- Gary – look at classification accuracy of proficiency and scale score – can do propensity score matching to look at types of students, look at impact of new speaking system – scoring methods – locally scored – construct irrelevant variance</p>		<ol style="list-style-type: none"> 1. In report - Need to draw out differences – 4 ways they are different – explicitly show effects of differences 2. Policy implications of new and better speaking scale. 							
A2.4	Speaking - reactions to demo	<p>Differences – currently face to face – here more scaffolding, model answer, different construct What is motivation for model answer – type of language and level of detail – may not currently expand and elaborate – T listening to answer, not language – cog lab, speak like nina speaks With one on one – teacher can answer questions – and scaffold – can you tell me more Evidence – classifications are very different there must be value added for disruption in amaos More meaningful generalizable performance you are obtaining from performance Validation strategy – what other kind of similar task can be administered to verify construct - expect higher correlation – only one construct – resources to express in each domain – not separate domain test s – hope correlations will go up Lyle - Several different types of language activities to demonstrate their capacity – each domain is not skills, but they are how students express themselves in language. Different methods Jamal - Issue of common core requirement more language – explain elaborate, discuss. Language in general is important not only academic language. Students proficient in general and academic Margo – interaction of multiple domains Domain is primary but is mediated through other domains. How to couch that – composite – one underlying construct – moving toward more authentic ways to capture language proficiency</p>		<ol style="list-style-type: none"> 1. Think about how to move towards more authentic ways to capture language proficiency. 							

A3.1	Writing - online and paper	<p>• State policy to decide whether students will take paper or online. Ability to demonstrate best language – know something about student – at end of first operational year – don't want to see differences if groups are different – some qualitative analyses on schools – whether they are different – instead match by proficiency level – student level – need to know why kids are assigned – school versus individual – is there evidence that student was assigned because they may do better.</p> <p>Once all schools offer cb, then only identified students will be accommodated. How to distinguish in linking, equating – identify students who would need accommodation.</p> <p>Cb may not be appropriate for all students.</p> <p>How to identify students with ability that requires paper based</p> <p>Are students performing differentially – there is an iep identifier in dataset – no iep, no accommodation</p> <p>In NH, decision not based on iep – teacher/team decision- parents resist having students identified.</p> <p>Accommodations offered with current paper – know what kids are getting and how will those be carried over to cb. Note – pb is not blanket accommodation</p> <p>Personal needs profile – appropriate accommodations provided seamlessly online</p> <p>Subcommittee looking at accommodations and how we ensure they are equitable</p> <p>Tim: will want to know what kinds of students take paper vs. online. If school decides this for all students rather than case-by-case (accommodations), how does this affect the study results?</p> <p>Dorry: This study can be done after the completion of the 1st operational year.</p> <p>Gary: students who don't have an IEP should not take the paper-based accommodation (so take out from the data set).</p> <p>Tim: students without IEP can still take something with an accommodation.</p> <p>Lyle: Consider what kinds of accommodations are offered for paper version now and what kind of accommodations are offered for paper future and online.</p> <p>• ACTION FOR WIDA: Pass studies done through the AAE SC to the TAC</p> <p>• Carol: Same raters vs. different raters: investigate whether the scoring result is the same</p> <p>• Lyle: keyboarding is construct relevant.</p> <p>• Jamal: add student survey question – Do you have a computer at home? (Too late to add to student survey for 2015 FT)</p> <p>• Lyle and Aki's research idea: Create a dummy variable to look at various groups: those who scored high but did not feel that they're better at keyboarding, those who scored low but feel that they're good at keyboarding, etc. Compare to their</p>		<p>1. Conduct study on students who take paper vs. online after completion of the 1st operational year - differentiate between students assigned based on individual differences vs classroom . (Sandra post-meeting thought: Add research questions about this in end of test student survey for operational test?)</p> <p>2. Think about how paper version is used whether it's for accommodations or other reasons. Consider what kinds of accommodations are offered for paper now and what kinds are offered for paper in the future.</p> <p>3. Bring studies done through the AAE SC to the TAC.</p> <p>4. Conduct additional study of same raters vs. different raters. Are results the same?</p>						
A3.1	Listening and Reading	<p>Two formats – not policy</p> <p>Comparability at scale score level</p> <p>Interpretations – how will we ensure interpretation based on two formats will be comparable</p> <p>Accommodation – should be alternative mode – schools not able to administer online also, not accommodation for all students- affordance rather than accommodation. Change the name of paper form from "Accommodated form" to just "paper form" or something else.</p> <p>Adjustments that do not alter construct – administration - from standards – don't alter the construct</p> <p>Calibrating listen and reading for online separately, then calibrate paper separately?</p> <p>Imposing parameters on paper test – are parameters similar? Note – not the same items</p> <p>Gary Require good sample for field test 350 is a little small</p> <p>Jamal May not be able to do item level analysis- items are different – DIF would not apply</p> <p>Can we include some external criteria – see if structures the same or difference – student performance on entire assessment – have operational access</p> <p>How many observed variable - 18-21 items for factor analysis – different items – claim standards are the same</p> <p>Create testlets based on contents – same test specs – distribution of math same as LASS,,,</p> <p>Expect loadings on LR PB to mirror CB</p> <p>exploratory or confirmatory- confirmatory to test variance</p> <p>Carol: same student s- within student analysis –2 parallel winsteps – look for comparable difficulty– rule out alternative explanation – means it is mode effect</p> <p>no linking items – common students</p> <p>ensure sample is good- across range of pls – particularly for confirmatory factor analysis</p> <p>there are a lot of items for small number of students</p> <p>could reduce to thematic folder?</p> <p>May not have enough students for confirmatory</p> <p>Number of factors depends on tier</p> <p>Scores dichotomous – do tetracorelation, use mplus.</p> <p>Alternatives to FA</p> <p>Scores should have same factor structure- same subtraits of domain- should have what we designed – unidimensional</p> <p>Comparability is a multifacet option</p>		<p>1. Topic should be: how will we ensure interpretation based on two formats will be comparable</p> <p>2. Be careful on use of word accommodation – support</p> <p>3. there are a lot of items for small number of students. Could reduce to thematic folder?</p> <p>4. Comparability with external criteria –try correlation with other domains.</p> <p>Run Two parallel winsteps for cb and pb – should be able to plot PLS- correlation</p> <p>5. Run confirmatory factor analysis – tetracorelation, use mplus.</p>					Change the name of paper form from "Accommodated form" to just "paper form" or something else.	
A3.1	Speaking - PB vs. CB	<p>Is mode study worth doing- resources – look at biggest problems – lit review should show there is no issue</p> <p>Creating an operational mode that never existed</p> <p>Half paper based, half cb, one test – counter balanced. responses recorded</p> <p>Either CAL or WIDA score</p> <p>Rater condition study - paper – locally scored will be scores of record. Send responses to be centrally scored. Should not be interaction with administrator</p> <p>Training of local scorers – training materials are being developed – TAs will be trained- scorer rubric – simplified for distributed scoring – pilot study – try it out at pilot – to be used for centralized scorers. Who are raters for locally scored – policy – may change a lot- that's why we made scoring simpler</p> <p>Other types of studies during quality control– take matched students –pick on centrally scored and locally scored –post facto analysis-</p> <p>Raters – should meet standards – but what happens locally is beyond our control</p> <p>Hard to make judgments about spoken language online- hope local raters will use model answer. Note – model answer may be at too high a level – students try to reach same level</p> <p>Lyle - Maybe remove model response down the road – may be part of tutorial– students like model in cob lab – not intimidating. Do they know what model is all about? Extrapolation – generalizability – artificial context – as test becomes better known – take that part out – may be interfering</p> <p>Other TAC members like the model response, and think it replaces the scaffolding of TAs.</p> <p>No followup or rewording of prompt currently – model takes the place of this option</p> <p>Speaking happens in context of conversation – exploring more interactive things – students talk to Nina</p> <p>Responses are edited so raters only hear responses – will raters know which questions – yes – will have all resources they need.</p>		<p>1. How is decision made for comparability – do we need to have all aspects of comparability – will be same task? Yes- what happens if one is not met?</p> <p>2. Interpretation of performances is highest level- that claim – people will</p> <p>Have to weigh whether issue affect interpretations</p> <p>3. Look at overall Ratings and scores</p> <p>Concurrent operational– if we know interaction for current than can we demonstrate they are more comparable for new- does rater task interact with mode – effects of interaction – is it better on new.</p> <p>4. Conduct lit review to show different mode does not present problem</p> <p>5. Conduct rater condition study - locally score response and send same response to be centrally scored</p> <p>6. Post facto analysis?</p> <p>7. Study on model response? (potential student survey question?)</p>						

A3.2	Keyboarding	<p>Extrapolation of generalizability rather than meaningfulness – is construct to generalize to domain in which they are writing</p> <p>We have new comers – sife – need to give students opportunity to build kb ability</p> <p>Policy – how are we defining this construct?</p> <p>Outcome – difference in score kb vs hw</p> <p>Is hw gold standard_ no what is to be compared against – what is changing</p> <p>May make a big difference for individual students</p> <p>Reliability of differences will be very low – what does small difference mean vs large difference – how to interpret</p> <p>Regression probably won't tell much – more latent class problem - use difference scores as predictor rather than outcome</p> <p>Could students take longer online?</p> <p>Summary of survey for ecss broken by grade/form</p> <p>Number of classes required degrees of freedom – only 350 variables</p> <p>Do two separate analyses – self identified – correlation with computer delivered vs correlation of hw</p> <p>Guidance to teachers of what to consider – brand new students</p> <p>6 variables – how many factors – instead look at very low vs very high – match as dummy variables – two groups of students – look at mean criterion score on test – may be variance in demographics</p> <p>Is this worthwhile? Yes – but keep it simple</p> <p>Propose to wida- replicating scoring study</p> <p>Keyboarding is future of assessment</p> <p>Writing score is most significant predictor of content scores</p>		<p>1. Is KB relationship construct irrelevant?</p> <p>Policy – how are we defining this construct?</p> <p>2. The keyboard indicator is only a</p> <p>3. will student level variables make an impact</p> <p>4. Regression probably won't tell much – more latent class problem - use difference scores as predictor rather than outcome</p> <p>is developing indicator of keyboard readiness worthwhile?</p> <p>The TAC says Yes – but keep it simple</p> <p>5. Do two separate analyses – self identified – correlation with computer delivered vs correlation of hw</p> <p>6. Propose to WIDA- replicating scoring study- TAC supports this idea- could same student enter in data? May be sufficient to just type students' responses.</p> <p>– cite studies that have done this in different areas – juxtapose different findings</p>										
A4	Standards Setting	<p>Two related but different purposes for standard settings</p> <p>First standard setting – determine a score that represents an English language proficiency criterion. A floor level below which a state may not go</p> <p>States may want to go above it – meeting common definition is level below which no state will go.</p> <p>Can you have multiple points for different grades – analysis between content performance and language proficiency will have to be smooth – set all grades to same point – everything is scaled to fit that level, so bar will be set at same level for all grades</p> <p>What do we use as impact data for bookmarking – if scale is maintained – oe historical access data- can go beyond just using field test data – put cuts by grade level- there are different methodologies – multiple parallel panels External agency – rfp to reach job. Give them parameters – what are critical issues to be thought about – what evidence will you be showing –</p> <p>Previous standard settings (Tim) problems – went with middle of road – not as connected with proficiency as with political agenda. Should higher eds be on panel for 6th grade – will have lower expectations results can be skewed by panelists – what does it mean to be ELL in class without support – need to be extremely well defined.</p> <p>Do we make recommendations about appropriate methodologies – Jamal – standard methods - bookmark, body of work, mapit based on set of arguments -how independent should CAL be? Very independent for first standard setting, and CAL very involved for second standard setting</p>	<p>~include state perceptions of policies that each of the 35 states have adopted. Get anecdotal feedback and provide as a background discourse document. This Qualitative data can help support empirical information to smooth over each grade level.</p> <p>~Panelists would then feel that they are being heard and considered.</p> <p>~Share with an articulation panel</p> <p>~Any scale scores that is re-worked should be able to put in terms of the current. Consider in the 2nd Standards Setting study. Maybe the 2nd Standards Setting study is the articulation panel. Qualifications for the 2nd study is more stringent because they need to have had student experience.</p>	<p>1. WIDA will contract a 3rd party to propose a study by July 2015 (1. Single cut score and performance standard for reclassification eligibility. 2. Multiple cut scores to interpret scale scores per WIDA ELD Standards) Check with Gary if Research has done anything on this. WIDA will need to issue an RFP announcement soon to do this!</p> <p>2. Definition of common cut score - A score that represents an English language proficiency criterion</p> <p>3. what does it mean to be ELL in class without support – need to be extremely well defined? Need to be really careful- potentially very political and very sensitive</p> <p>4. what do we know about scores, and how successful they are – other measures – feed into first standard setting – what would that be?</p> <p>5. (TPS) Include state perceptions of policies that each of the 35 states have adopted. Get anecdotal feedback and provide as a background discourse document. This Qualitative data can help support empirical information to smooth over each grade</p>										
B1	Drift	<p>o Lyle's food for thought questions:</p> <p>How many of "those" do you need to address?</p> <p>Is ACCESS leading reclassification decisions?</p> <p>How does this affect comparability?</p> <p>How many stakeholders would be concerned about this? Ask for their feedback. Have we done the best we can?</p> <p>o Carol: overwhelmed by MT's report. It was trying to create best case scenario rather than year to year overtime analysis. 2/3 of rerated scores show increases – results do not seem to support conclusion that there is no drift</p>		<p>1. Get information from states about whether it impacts students – are more students exiting?</p>										
B2	ATR	<p>Carol- Much better – clearer, easier to read</p> <p>Lyle – want focus to turn to full AUA</p> <p>• For ACCESS 2.0, we'll revisit the design of the ATR.</p>		<p>1. Get feedback from states and other users (TPS?)</p> <p>2. Discuss how to turn towards AUA in the ATR</p> <p>3. Re-design ATR for 2.0</p>				3. Re-design ATR for 2.0						
C1	Alt ACCESS	<p>Small differences between equating and population give TAC confidence in small sample in field test</p> <p>TAC is ok with leaving series 102 alone and update series 103. No new items</p>		<p>1. Leave series 102 alone and update 103.</p>										
	Early Years	<p>Once you have this early years language assessment going, would it have potential for ACCESS? – could it be used instead of using listening and reading scores to indicate where to start on ACCESS online speaking test? –</p> <p>For students designating at very beginning – there is a home background, other information beyond screener improve accuracy of screener- home language survey – access is only one piece of information – multifaceted.</p> <p>• Lyle: No need for separate AUAs for each component of the assessment</p> <p>• Lyle: this is an excellent start. Will share comments with Alicia</p> <p>• Lyle's group has developed a software to help create AUA</p>	<p>under what circumstances would you choose which type of the screener? Need to be clear.</p>	<p>1. Need claim 3 for each screener</p> <p>Maybe claim 4.</p> <p>2. Under what circumstances would you choose which type of the screener? Need to be clear.</p>										
	Other TAC meeting-related suggestions	<p>Margo- Communication – consider condensing and sharing the readaheads with states to demonstrate how thoroughly we have thought through issues and types of analysis – white paper? – one pager per topic?</p> <p>Dorry - Jenn start to build aua for field test</p>	<p>~Advice to WIDA and CAL: Avoid informational discussion. At least present 3 questions per topic to get to the heart of the matter.</p> <p>~Allow TPS SC members to speak at TAC meeting</p> <p>~Should DRC participate? What is their role in the TAC and TPS?</p> <p>~Request DRC to provide weekly reports showing evidence of quality control</p>	<p>1. Margo- Communication – consider condensing and sharing the readaheads with states to demonstrate how thoroughly we have thought through issues and types of analysis – white paper? – one pager per topic?</p> <p>2. Jenn start to build AUA for Field Test</p> <p>3. (TPS) Next TAC meeting, avoid informational discussion and allow TPS members to speak.</p> <p>4. Determine whether DRC should participate and to what degree?</p> <p>5. (TPS) Provide list of reports that WIDA is already requesting of DRC and provide to TAC and TPS for their sign-off (reports show evidence of QC)</p>										
TPS	TPS Purpose and Responsibilities	<p>• Purpose: Get more state input in technical decisions</p> <p>• Responsibilities:</p> <p>~Robert Lee: Need to get heads together about measuring assessments. Work with Gary Cook on this</p> <p>~Report back to the EC. Simplify, translate, and amplify to the EC. When we propose ideas to EC, it is after consultation with TPS members</p> <p>~Pass information from respective states to WIDA and CAL</p> <p>~Provide feedback on read-ahead materials</p>		<p>1. Write Mission Statement and send to members</p> <p>2. Vendor QC - consider inviting a SC member</p> <p>3. CAL QC review - consider inviting a SC member</p>										

	TPS Structure	<p>~Ideal number of TPS members is 6-9 representatives</p> <p>~want to have TPS be able to speak up at TAC</p> <p>~How to report out to the EC?</p> <ul style="list-style-type: none"> ~One person in this group will report to the EC in addition to some kind of read-ahead. ~Steve proposes a rotating representative ~Can join by phone or in person ~Half an hour report <p>• Name of the SC</p> <ul style="list-style-type: none"> ~Psychometric and Policy Advisory SC (PPASC) ~Psychometric Advisory SC (PASC) ~Assessment Design and Development Advisory SC (ADDA) ~Measurement Advisory SC (MASC) ~Assessment Policy Advisory SC (APASC) ~Validation Advisory SC <p>• Timeline of Meetings:</p> <ul style="list-style-type: none"> ~ Have Quarterly Meetings (counting the 2 TAC meetings) ~other 2 meetings are in the form of webinars (Two 1.5 hr webinars) to discuss issues (1st 3 bullets of the responsibilities slide) ~Subcommittee help come up with some agenda topics ~Standing Topic: Reviewing Equating Procedures ~Late Feb 2015 potentially ~Face-to-Face in the Fall TAC (Nov. or Dec.) 		<p>1. Ask SC to vote on name of the group when we send out the Mission Statement</p> <p>2. Send Doodle poll to set date for next meeting (late Feb. 2015 potentially) and provide info on next TAC dates.</p> <p>3. Invite a TPS member to attend EC via call?</p>										
	Next TPS Meeting Agenda	<p>~1st Webinar discuss objectives for the year.</p> <p>~Think about what info you would like from the new Test Development vendor (Dave Chayer)</p> <p>~Standing Topic: Reviewing Equating Procedures</p> <p>~Research on growth model (Gary's work) and providing feedback (priority as expressed by Steve and Robert)</p> <p>~Take in front of the TAC and the SC measures of growth in order to present full recommendation to the EC. Urge Gary to pilot student growth percentiles in 2015. Consortium-level growth percentiles.</p> <p>~QC – explore 3rd party QC procedure</p> <ul style="list-style-type: none"> ~MI has AES conduct the QC. Examine if the delta is big enough ~WIDA would like to issue an RFP. The SC can potentially help! ~Who is the arbiter? This committee. <p>~Quality control plans for automated scoring and hand-scoring. Help solidify the language from the DRC RFP.</p> <p>~With embedded field test, we will be able to use data from the field test to do pre-equating. Discuss what near future equating will look like</p> <p>~Do we want a larger equating sample? How does this work with the reporting window?</p> <ul style="list-style-type: none"> ~See what criteria is needed to determine sample <p>~MA and MI both have done this: Consider post-equating and then used that as parameters for pre-equating. Write up and submit to the TAC. This can mitigate drift. Field Test data is set against the post-equating data rather than pre-equating data.</p> <ul style="list-style-type: none"> ~Acknowledging a better estimate ~Work on a proposal together in the next SC Quarterly Meeting call. <p>~Review the blueprint design of the embedded Field Test (test map, refreshment schedule, test specifications) – Also read the Test and Item Design Plan.</p> <ul style="list-style-type: none"> ~Randomly selected students receive embedded field test item <p>~Innovative item types: hot spots, drag and drop, etc. What do these mean for placement?</p> <p>~Alternate ACCESS – will not be on the same scale. Discuss Proficiency Level Interpretation. Work on new version of Alternate ACCESS.</p>		<p>1. Shu Jing: CAL can send PPT presentation from 3 years ago that explains the process of equating. Keep equating sample as stable as possible.</p> <p>2. Consider post-equating and then used that as parameters for pre-equating. Write up and submit to the TAC. This can mitigate drift. Field Test data is set against the post-equating data rather than pre-equating data.</p> <p>3. Carsten send Steve (rest of SC) document about current Proficiency Level.</p>										