

## Supplemental Materials

### Detailed descriptive statistics of the children's proficiency scores

As mentioned in the main paper, out of the 49 children with proficiency data: 14 children had equal comprehension proficiency in both languages (8 French–English, 6 Spanish–English); 17 children were more proficient in English comprehension than French/Spanish (12 French–English, 5 Spanish–English); and 18 children were more proficient in French/Spanish than in English (10 French–English, 8 Spanish–English). Proficiency data was missing from 1 French–English and 3 Spanish–English children. Table S1 lists the detailed descriptive statistics of the proficiency scores.

**Table S1.** *Language proficiency scores across the two groups of bilingual children. Proficiency data were missing for 1 French–English and 3 Spanish–English bilingual children.*

	French–English			Spanish–English		
	N (%)	Mean (SD)	Range	N (%)	Mean (SD)	Range
<b>Overall</b>	30 (97%)			19 (86%)		
English		8.60 (1.71)	3 – 10		8.32 (2.16)	3 – 10
French/Spanish		8.23 (2.14)	3 – 10		9.16 (0.90)	8 – 10
<b>Equal proficiency</b>	8 (26%)	9.25 (0.89)	8 – 10	6 (27%)	9.67 (0.82)	8 – 10
<b>More proficient in English</b>	12 (39%)			5 (23%)		
English		9.42 (0.90)	8 – 10		9.60 (0.55)	9 – 10
French/Spanish		6.42 (2.27)	3 – 9		8.20 (0.45)	8 – 9
<b>More proficient in French/Spanish</b>	10 (32%)			8 (36%)		
English		7.10 (2.02)	3 – 9		6.50 (2.20)	3 – 9
French/Spanish		9.60 (0.70)	8 – 10		9.38 (0.74)	8 – 10

A linear regression model with proficiency score as the dependent variable was run to compare proficiency scores between the French–English and Spanish–English bilinguals. In the linear regression model, fixed effects included language community (French–English vs. Spanish–English) and language (English vs. French/Spanish), as well as their interaction<sup>1</sup>:

$$\text{proficiency} \sim \text{language} * \text{lang\_community}$$

<sup>1</sup> We tried running a linear mixed-effects model with a random intercept for participants; however, the model returned a singular fit.

We found no significant effect of language, language community, nor their interaction ( $ps > 0.109$ ). Therefore, there was no significant difference between the French–English and Spanish–English children in terms of their proficiency in English or their proficiency in the other language (i.e., French for the French–English children and Spanish for the Spanish–English children).

### **Exploratory analysis of the main paper**

In the main paper, we reported the preregistered analyses on accuracy — our primary dependent variable in determining bilingual children’s word learning in touching the labeled target object on each test trial. In this supplemental materials, we report two preregistered exploratory analyses: (1) looking at the effect of age, (2) using response time as the dependent variable, and an additional analysis on the effect of language proficiency.

#### ***Effect of age***

Our accuracy analysis revealed that bilingual children in both communities successfully learned the novel words regardless of the language switching patterns used during the learning blocks. As previous research has shown that children’s ability to learn words may improve with age (e.g., Read et al., 2021; Scaff et al., 2022), in our preregistration we also expected older children to show a greater accuracy than younger children. Therefore, we additionally compared models with and without age as a predictor variable. Bilingual children’s age in months was scaled and centered for ease of interpretation. The final model specification was:

$$\text{accuracy} \sim \text{condition} * \text{lang\_community} * \text{age\_in\_months} + (1 + \text{condition} | \text{participant}) + (1 | \text{item})$$

When added to the model, neither the main effect nor interactions with age were significant (all  $ps > 0.302$ ; see Table S2 for the coefficient estimates from this model). Moreover, a model comparison with and without age as a variable indicated no significant improvement in model fit,  $\chi^2(4) = 2.26$ ,  $p = 0.687$ . Overall, the pattern in this model was consistent with the main accuracy analysis reported in the paper where there was no significant difference in terms of

condition or language community, suggesting that bilingual children successfully learned the novel words regardless of age.

**Table S2.** Coefficient estimates from the logistic mixed-effects models predicting accuracy in the test blocks with *age\_in\_months* as an additional fixed effect.

	Estimate	SE	z	p
Intercept	1.24	0.22	5.61	<.001
condition	0.12	0.31	0.39	0.700
lang_community	0.24	0.38	0.62	0.537
age_in_months	0.11	0.18	0.60	0.546
condition * lang_community	0.02	0.59	0.03	0.973
condition * age_in_months	0.09	0.29	0.30	0.766
lang_community * age_in_months	-0.26	0.36	-0.71	0.479
condition * lang_community * age_in_months	0.61	0.59	1.03	0.302

### **Response time**

In addition to accuracy, we explored response time on each correctly-responded test trial as a dependent variable; a total of 853 correct trials were included in this analysis. The decision to use response time as an additional measure was driven by a possibility that this measure might be able to better capture more individual variability in terms of the speed of children's lexical comprehension performance, since children were generally very accurate in the task. Moreover, response time may also be more sensitive to age effects, as children answer more quickly as they get older and gain more expertise with language (Scaff et al., 2022). On average, French–English bilingual children had a mean response time of 1997ms in the *immediate-translation* condition ( $SD = 1064.20$ ; range = 451.86 – 5130.50) and 1788ms in the *one-language-at-a-time* condition ( $SD = 880.09$ ; range = 339.33 – 4323.33). On the other hand, Spanish–English bilingual children had a mean response time of 2620ms in the *immediate-translation* condition ( $SD = 1353.25$ ; range = 764.33 – 6166.50) and 2303ms in the *one-language-at-a-time* condition ( $SD = 1455.58$ ; range = 457.25 – 5550.60).

A linear mixed-effects model was run on response time. To correct for issues of non-normality, raw response time was log-transformed. Condition, language community, and age, as

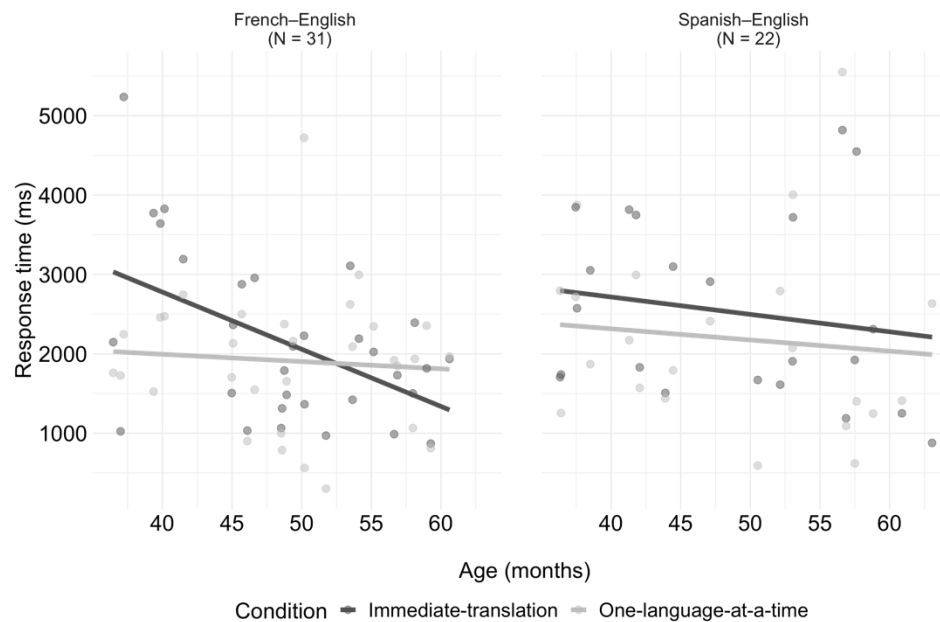
well as their interactions, were entered as fixed effects in the model; a random slope of condition by participant and a random intercept of stimulus item were also entered:

$$\log\_rt \sim \text{condition} * \text{lang\_community} * \text{age\_in\_month} + (1 + \text{condition} | \text{participant}) + (1 | \text{item})$$

The coefficient estimates from this model are shown in Table S3 and Figure S1 visualizes this model. The only terms that approached significance was the three-way interaction. Following up on this three-way interaction, separate linear mixed-effects analyses were run for the French–English bilinguals and the Spanish–English bilinguals. The models revealed that the effect of age approached significance in the *immediate-translation* condition for the French–English bilinguals (Estimate = -0.19, SE = 0.10,  $t = -1.76$ ,  $p = 0.088$ ), suggesting that the reaction time for French–English bilinguals decreased significantly in the *immediate-translation* condition across age. However, the effect of age was not significant for the Spanish–English bilinguals (Estimate = 0.02, SE = 0.11,  $t = 0.13$ ,  $p = 0.900$ ). Moreover, we found no significant interaction across age between the two conditions in either group of bilinguals. Overall, similar to the patterns reported in the first set of analysis on accuracy, we did not observe any significant difference in terms of condition or language community, nor in their interaction. Therefore, consistent with the accuracy analysis, bilingual children in both communities performed similarly in word learning across both the *immediate-translation* and *one-language-at-a-time* conditions.

**Table S3.** Coefficient estimates from the linear mixed-effects model predicting log-transformed response time in the test blocks.

	Estimate	SE	$t$	$p$
Intercept	7.22	0.08	94.60	<.001
condition	-0.10	0.08	-1.25	0.220
lang_community	0.16	0.15	1.13	0.265
age_in_months	-0.09	0.07	-1.33	0.189
condition * lang_community	-0.12	0.16	-0.78	0.443
condition * age_in_months	0.01	0.08	0.13	0.901
lang_community * age_in_months	0.07	0.14	0.51	0.616
condition * lang_community * age_in_months	-0.28	0.16	-1.79	0.080



**Figure S1.** Response time by condition, language community, and age in the test block. Individual dots plot the data from each individual participant.

### **Language proficiency**

In addition to the preregistered analyses, we explored the effect of language proficiency on bilingual children's learning of the novel cross-language words. Previous research has revealed mixed evidence as to whether language proficiency interacts with bilingual children's word learning ability during different types of bilingual book reading sessions (e.g., Brouillard et al., 2022; Read et al., 2021). Therefore, it is plausible that language proficiency may have an effect on how bilingual children learn from different language switching patterns.

Building upon the logistic mixed-effects model used in the main accuracy analysis, we added a variable of proficiency score to the model. This created a continuous variable with the caregiver-reported proficiency rating in each language, which we used to predict children's performance on trials in that same language. In this analysis, we excluded data from one French-English and three Spanish-English children participants who were missing their proficiency information; data from 30 French-English and 19 Spanish-English children remained in the analysis. The initial specification was:

$$\text{accuracy} \sim \text{condition} * \text{lang\_community} * \text{proficiency} + (1 + \text{condition}|\text{participant}) + (1|\text{item})$$

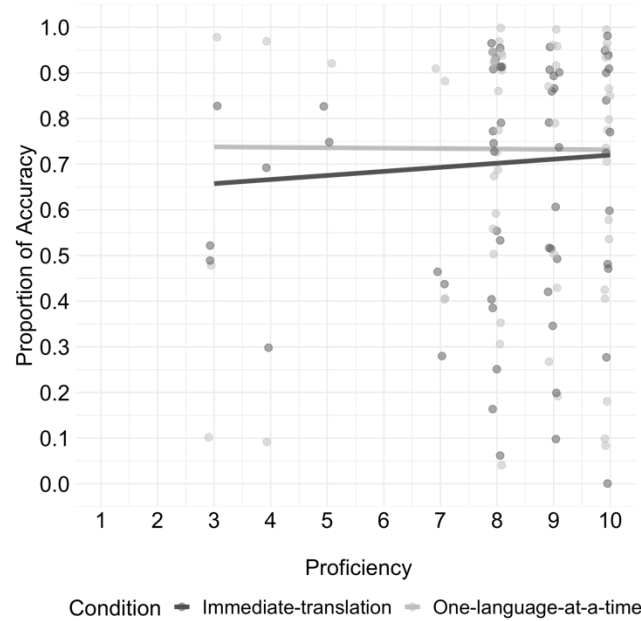
However, as the initial model did not converge, we first removed the random slope for condition and the random intercept for item. Moreover, since we did not find any significant difference between the two communities in the main accuracy analysis, we performed a model comparison between the model with language communities and the one without. This comparison also indicated no significant improvement in model fit,  $\chi^2(4) = 4.74$ ,  $p = 0.315$ , so we further pruned the effect of language community from the model. Therefore, the final model was:

$$\text{accuracy} \sim \text{condition} * \text{proficiency} + (1|\text{participant})$$

The coefficient estimates from this model are shown in Table S4 and Figure S2 visualizes this model. Similar to the main analysis, no significant difference between conditions was found. Moreover, we did not observe any significant effect of language proficiency. As can be seen in Figure S2, the level of proficiency did not hugely impact children's accuracy in our experiment. Note that we also ran another model including all the children who participated in our experiment (i.e., including those who did not initially meet our language proficiency criteria; please refer to the section below for the detailed statistics).

**Table S4.** Coefficient estimates from the logistic mixed-effects model predicting accuracy in the test blocks with language proficiency scores.

	Estimate	SE	z	p
Intercept	0.53	0.47	1.14	0.255
condition	0.65	0.69	0.94	0.347
proficiency	0.08	0.05	1.48	0.138
condition * proficiency	-0.06	0.08	-0.80	0.423



**Figure S2.** *Proportion of accuracy by condition and language proficiency in the test blocks. Individual dots plot the data from each individual participant.*

### Supplemental analysis using the preregistered exclusion criteria

Analyses reported in the main paper deviated from the preregistered language exclusion criteria, as it resulted in exclusion of a higher than anticipated number of children and thus led to a smaller sample size and decreased statistical power. For transparency, this supplemental material reports the analyses using the more stringent preregistered exclusion criteria.

Following the preregistered language exclusion criteria, a total of 10 French–English and 12 Spanish–English children were excluded. When additional exclusion criteria were applied (see participants section in the main paper), the remaining sample consisted of 22 French–English children (13 girls; Mean age = 4.04 years,  $SD = 0.56$ , range = 3.04 – 4.94) and 14 Spanish–English children (7 girls; Mean age = 4.21 years,  $SD = 0.75$ , range = 3.03 – 5.26).

Language proficiency information was missing for 1 French–English and 3 Spanish–English children; 13 children had equal comprehension proficiency in both languages (7 French–English, 6 Spanish–English); 6 French–English children were more proficient in English comprehension than French/Spanish; and 13 children were more proficient in French/Spanish

than in English (8 French–English, 5 Spanish–English). Table S5 contains descriptive statistics of the proficiency scores. We also ran a linear regression model with proficiency score as the dependent variable to compare proficiency scores between the French–English and Spanish–English bilinguals<sup>2</sup>. There was no significant effect of language, language community, nor their interaction ( $ps > 0.301$ ). Similar to the sample reported in the main paper, there was no significant difference between the French–English and Spanish–English children in terms of their proficiency in English as well as their proficiency in the other language (i.e., French for the French–English children and Spanish for the Spanish–English children). In this more restricted data set, 68% of the mothers in Montreal and 64% of the mothers in New Jersey had completed a university degree or higher.

**Table S5.** *Language proficiency scores across the two groups of bilingual children whose proficiency scores were available.*

	French–English			Spanish–English		
	N (%)	Mean (SD)	Range	N (%)	Mean (SD)	Range
<b>Overall</b>	21 (95%)			11 (79%)		
English		9.00 (1.10)	7 – 10		8.91 (1.14)	7 – 10
French/Spanish		9.19 (0.93)	8 – 10		9.64 (0.67)	8 – 10
<b>Equal proficiency</b>	7 (32%)	9.29 (0.95)	8 – 10	6 (43%)	9.67 (0.82)	8 – 10
<b>More proficient in English</b>	6 (27%)			0 (0%)		
English		10.00 (0.00)	10		—	—
French/Spanish		8.33 (0.52)	8 – 9		—	—
<b>More proficient in French/Spanish</b>	8 (36%)			5 (36%)		
English		8.00 (0.76)	7 – 9		8.00 (0.71)	7 – 9
French/Spanish		9.75 (0.71)	8 – 10		9.60 (0.55)	9 – 10

In the following analyses, we followed the same procedure as in the main paper, where we first explored the preregistered analyses with proportion of accuracy as the dependent variable and then performed exploratory analyses: (1) on the effect of age, (2) with response time as the dependent variable, and (3) the effect of proficiency.

<sup>2</sup> In the linear regression model with proficiency scores as the dependent variable, fixed effects included language community (French–English vs. Spanish–English) and language (English vs. French/Spanish), as well as their interaction.



### **Accuracy**

**Familiar word block.** French–English bilingual children showed a mean accuracy of 0.98 in the familiar English-word trials ( $SD = 0.05$ ; range = 0.83 – 1.00) and 0.98 in the familiar French-word trials ( $SD = 0.05$ ; range = 0.83 – 1.00). Meanwhile, Spanish–English bilingual children showed a mean accuracy of 1 in the familiar English-word trials ( $SD = 0.00$ ) and 1 in the familiar Spanish-word trials ( $SD = 0.00$ ). As children’s performance was almost at ceiling with little variance, it was not possible to fit the logistic mixed-effects model to compare performance across the two communities. On the other hand, this near-ceiling accuracy also suggests that our preregistered exclusion criteria could be too stringent such that only children with nearly perfect accuracy were included in the analyses.

**Test blocks.** On average, French–English bilingual children showed a mean accuracy of 0.72 in the *immediate-translation* condition ( $SD = 0.25$ ; range = 0.17 – 1.00) and 0.68 in the *one-language-at-a-time* condition ( $SD = 0.27$ ; range = 0.00 – 1.00). On the other hand, Spanish–English bilingual children showed a mean accuracy of 0.76 in the *immediate-translation* condition ( $SD = 0.23$ ; range = 0.25 – 1.00) and 0.70 in the *one-language-at-a-time* condition ( $SD = 0.24$ ; range = 0.33 – 1.00). Separate one-sample t-tests were run on the proportion of accuracy in each condition per community, and confirmed that children from both communities learned the novel words in each condition significantly above the at-chance level of 0.50 ( $ps < .01$ )<sup>3</sup>.

Following the main paper, we ran a logistic mixed-effects model on the proportion of accuracy, with condition and language community as fixed effects, and a random slope of condition by participants and random intercept of item:

---

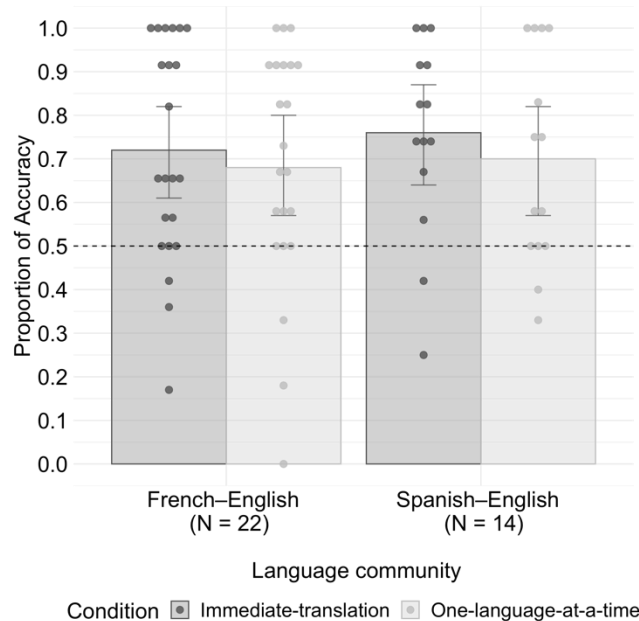
<sup>3</sup> For the French–English bilinguals, they performed significantly above the at-chance level in the *immediate-translation* condition ( $t(21) = 4.01$ ,  $p < .001$ ) and the *one-language-at-a-time* condition ( $t(21) = 3.14$ ,  $p = .002$ ). Likewise, the Spanish–English bilinguals performed significantly above the at-chance level in the *immediate-translation* condition ( $t(13) = 4.24$ ,  $p < .001$ ) and the *one-language-at-a-time* condition ( $t(13) = 3.03$ ,  $p = .005$ ).

$$\text{accuracy} \sim \text{condition} * \text{lang\_community} + (1 + \text{condition} | \text{participant}) + (1 | \text{item})$$

The coefficient estimates from this model are shown in Table S6, and Figure S3 visualizes this model. Consistent with the patterns found in the main paper, the model also did not reveal any significant difference in terms of condition or language community, and the interaction between condition and language community was also not significant. Therefore, similar to the main paper, bilingual children in both communities showed strong evidence of word learning in both language switching conditions.

**Table S6.** Coefficient estimates from the logistic mixed-effects models predicting accuracy in the test phase.

	Estimate	SE	z	p
Intercept	1.17	0.26	4.46	<.001
condition	-0.25	0.38	-0.65	0.513
lang_community	0.30	0.47	0.64	0.523
condition * lang_community	0.07	0.73	0.09	0.929



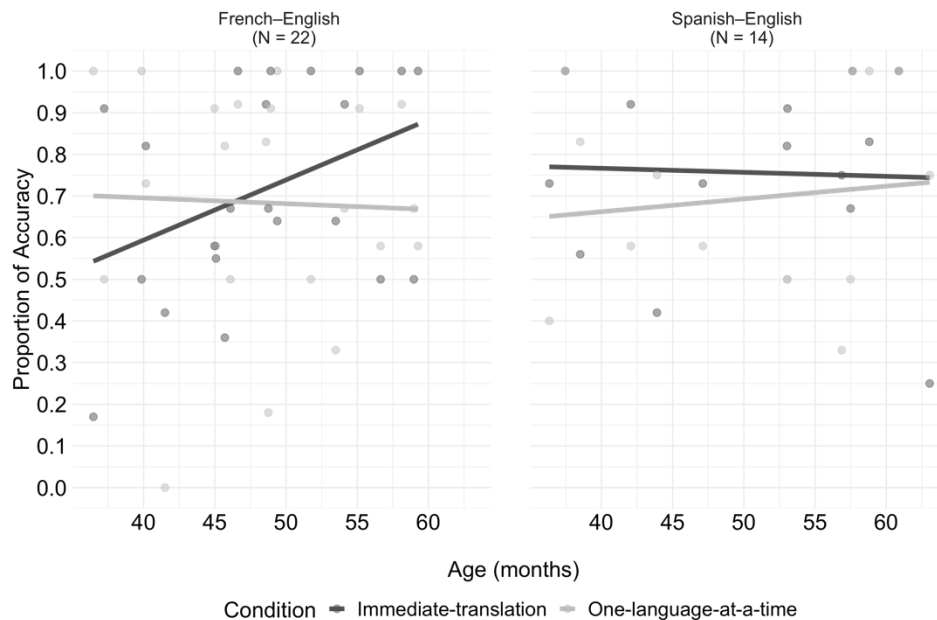
**Figure S3.** Average proportion of accuracy by condition and language community in the test blocks. Dots plot the data from each individual participant. Error bars indicate 95% confidence intervals, and the black dashed line represents the at-chance accuracy level of 0.50.

### Exploratory analysis

**Effect of age.** As an exploratory analysis, we also ran a logistic mixed-effects model with age as a fixed effect. Similar to the analysis on the sample reported in the main paper, model comparison with the model without age showed that adding age did not significantly improve the model ( $\chi^2(4) = 5.55, p = 0.236$ ). Moreover, we had to prune the random slope for condition from the model since it would not converge, the final model specification was:

$$\text{accuracy} \sim \text{condition} * \text{lang\_community} * \text{age\_in\_months} + (1|\text{participant}) + (1|\text{item})$$

The coefficient estimates from this model are shown in Table S7 and Figure S4 visualizes this model. The model did not reveal a significant effect of age, but there was a significant three-way condition \* language community \* age interaction. This was due to the crossover interaction in the French–English bilingual children, where their performance significantly improved on the *immediate-translation* condition with age. In contrast, the Spanish–English children showed similar performance in the two conditions and across age.



**Figure S4.** Proportion of accuracy by condition, language community, and age in the test phase. Individual dots plot the data from each individual participant.

**Table S7.** Coefficient estimates from the logistic mixed-effects models predicting accuracy in the test phase with *age\_in\_months* as an additional fixed effect.

	Estimate	SE	z	p
Intercept	1.06	0.24	4.49	<.001
condition	-0.33	0.18	-1.85	0.065
lang_community	0.24	0.42	0.56	0.574
age_in_months	0.20	0.19	1.05	0.295
condition * lang_community	-0.05	0.35	-0.15	0.878
condition * age_in_months	-0.30	0.18	-1.67	0.095
lang_community * age_in_months	-0.29	0.38	-0.75	0.455
condition * lang_community * age_in_months	0.98	0.36	2.73	<.01

**Response time.** We also included response time on each correctly-responded test trial as a dependent variable. There were a total of 573 trials in this analysis. On average, French–English bilingual children had a mean response time of 2101ms in the *immediate-translation* condition ( $SD = 1101.21$ ; range = 451.86 – 5130.50) and 1803ms in the *one-language-at-a-time* condition ( $SD = 765.06$ ; range = 339.33 – 3101.00). On the other hand, Spanish–English bilingual children had a mean response time of 2410ms in the *immediate-translation* condition ( $SD = 1108.38$ ; range = 764.33 – 4549.00) and 2033ms in the *one-language-at-a-time* condition ( $SD = 1371.97$ ; range = 457.25 – 5525.80).

We ran a linear mixed-effects model with condition, language community, age, as well as their interactions, entered as fixed effects. A random slope of condition by participants and random intercepts of item were also entered. To correct issues of non-normality in response time, raw response time was log-transformed. The final model specification was:

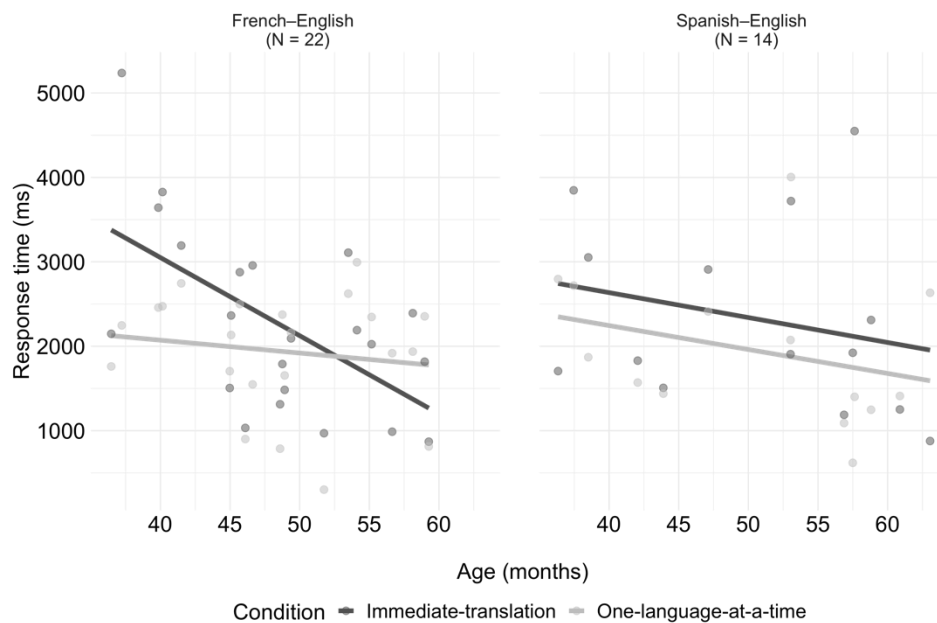
$$\log\_rt \sim \text{condition} * \text{lang\_community} * \text{age\_in\_months} + (1 + \text{condition} | \text{participant}) + (1 | \text{item})$$

The coefficient estimates from this model are shown in Table S8 and Figure S5 visualizes this model. Overall, visualization of the model suggests that, across both conditions and language communities, older children generally responded faster than younger children. However, the model did not reveal any significant effects or interactions of condition, language community, or age. Therefore, consistent to the patterns reported for the sample included in the main paper, in

terms of response time, bilingual children in both communities performed similarly in word learning across both the *immediate-translation* and *one-language-at-a-time* conditions.

**Table S8.** Coefficient estimates from the linear mixed-effects model predicting log-transformed response time in the test phase.

	Estimate	SE	<i>t</i>	<i>p</i>
Intercept	7.18	0.09	84.80	<.001
condition	-0.14	0.09	-1.53	0.137
lang_community	0.07	0.17	0.39	0.696
age_in_months	-0.13	0.08	-1.57	0.126
condition * lang_community	-0.17	0.18	-0.95	0.349
condition * age_in_months	0.04	0.09	0.40	0.690
lang_community * age_in_months	0.08	0.16	0.47	0.644
condition * lang_community * age_in_months	-0.19	0.18	-1.06	0.298



**Figure S5.** Response time by condition, language community, and age in the test phase. Individual dots plot the data from each individual participant.

**Language proficiency in all participants.** For a more inclusive analysis on the effect of proficiency, this analysis also included data from children who were previously eliminated for not fulfilling the language criteria — either the criteria reported in the main paper or the more stringent criteria reported in this document. In total, there were 30 French–English children and 22 Spanish–English children who were born full term and without any reported language

problems or participated without any problem. Among this group of children, the French–English children had a mean proficiency score of 8.42 ( $SD = 1.93$ , range = 3 – 10) and the Spanish–English children had a mean proficiency score of 8.34 ( $SD = 2.07$ , range = 2 – 10). To explore the effect of proficiency, we added a variable of proficiency score to the logistic mixed-effects model. The initial model specification was:

$$\text{accuracy} \sim \text{condition} * \text{lang\_community} * \text{proficiency} + (1 + \text{condition} | \text{participant}) + (1 | \text{item})$$

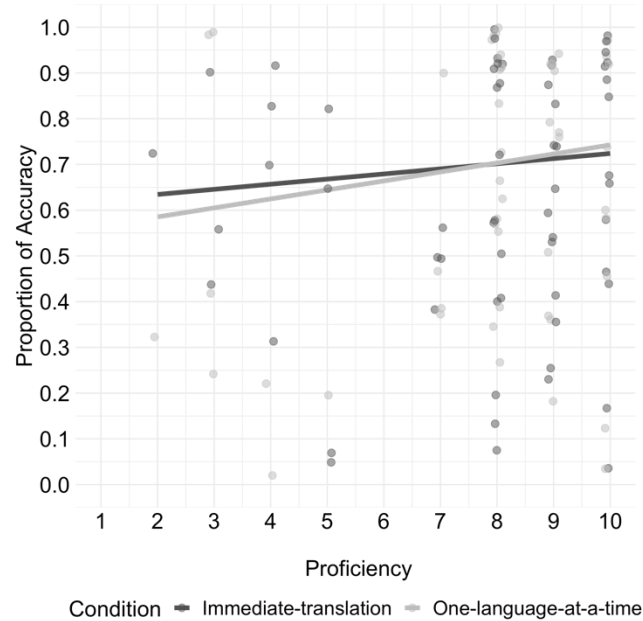
However, as the initial model could not converge, we removed the random slope for condition and the random intercept for stimulus item. Moreover, since we did not find any significant difference between the two communities in the main accuracy analysis, we further pruned the effect of language community from the model. Note that model comparison between the model with language community and the one without also indicated no significant improvement in model fit,  $\chi^2(4) = 7.44$ ,  $p = 0.12$ . Therefore, the final model was:

$$\text{accuracy} \sim \text{condition} * \text{proficiency} + (1 | \text{participant})$$

The coefficient estimates from this model are shown in Table S9 and Figure S6 visualizes this model. The model revealed a significant effect of proficiency, suggesting that bilingual children were overall more accurate for trials in which they had a higher level of language proficiency in the trial language. In contrast to the patterns reported for the sample included in the main paper, this model here revealed that the effect of proficiency may be more evident when a wider range of language proficiency level was included in the analysis. In other words, the sample reported in the main paper could have possibly focused on children who were relatively proficient. Yet, the lack of a significant interaction points to the possibility that proficiency level did not affect children's accuracy in learning words under different language-switching patterns.

**Table S9.** Coefficient estimates from the logistic mixed-effects model predicting accuracy in the test phase with proficiency scores among all the children.

	Estimate	SE	z	p
Intercept	0.11	0.42	0.26	0.794
condition	-0.43	0.60	-0.72	0.473
proficiency	0.12	0.05	2.63	<.01
condition * proficiency	0.05	0.07	0.70	0.485



**Figure S6.** Proportion of accuracy of all children by condition and proficiency in the test blocks. Individual dots plots the data from each individual participant.