

Study details	Population	Interventions	Study outcomes	Appraisal and Funding
<p>added value of these seems to be limited.</p> <p>Study dates Participants started gender-affirming therapy between 2001 and 2011</p>	<p><15 years and the old transfemales group ≥15 years.</p>		<p>Statistically significant increase (p≤0.05)</p> <p>Transmales (bone age <14 years), change from starting gender-affirming hormones to 24 months follow-up. Median (range), g/m³</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: 0.23 (0.19 to 0.28) • 24-months: 0.25 (0.22 to 0.28) • Statistically significant increase (p≤0.01) <p>z-score (range)</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: -0.84 (-2.2 to 0.87) • 24-months: -0.15 (-1.38 to 0.94) <p>Statistically significant increase (p≤0.01)</p> <p>Transmales (bone age ≥14 years), change from starting gender-affirming hormones to 24 months follow-up. Median (range), g/m³</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: 0.24 (0.20 to 0.28) • 24-months: 0.25 (0.21 to 0.30) • Statistically significant increase (p≤0.01) <p>z-score (range)</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: -0.29 (-2.28 to 0.90) • 24-months: -0.06 (-1.75 to 1.61) <p>Statistically significant increase (p≤0.01)</p> <p>Bone density: femoral neck</p> <p>Femoral neck BMAD</p> <p>Transfemales (bone age <15 years), change from starting gender-affirming hormones to 24 months follow-up. Median (range), g/m³</p>	

Study details	Population	Interventions	Study outcomes	Appraisal and Funding
			<ul style="list-style-type: none"> • Start of gender-affirming hormones: 0.27 (0.20 to 0.33) • 24-months: 0.27 (0.20 to 0.36) • No statistically significant change z-score (range) • Start of gender-affirming hormones: -1.32 (-3.39 to 0.21) • 24-months: -1.30 (-3.51 to 0.92) • No statistically significant change <p>Transfemales (bone age ≥ 15 years), change from starting gender-affirming hormones to 24 months follow-up.</p> <p>Median (range), g/m³</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: 0.30 (0.26 to 0.34) • 24-months: 0.29 (0.24 to 0.38) • No statistically significant change z-score (range) • Start of gender-affirming hormones: -0.36 (-1.50 to 0.46) • 24-months: -0.56 (-2.17 to 1.29) • No statistically significant change <p>Transmales (bone age <14 years), change from starting gender-affirming hormones to 24 months follow-up.</p> <p>Median (range), g/m³</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: 0.30 (0.22 to 0.35) • 24-months: 0.33 (0.23 to 0.37) • Statistically significant increase (p≤ 0.01) z-score (range) • Start of gender-affirming hormones: -0.37 (-2.28 to 0.47) • 24-months: -0.37 (-2.03 to 0.85) 	

Study details	Population	Interventions	Study outcomes	Appraisal and Funding
			<ul style="list-style-type: none"> • Statistically significant increase ($p \leq 0.01$) <p>Transmales (bone age ≥ 14 years), change from starting gender-affirming hormones to 24 months follow-up.</p> <ul style="list-style-type: none"> • Start of gender-affirming hormones: 0.30 (0.23 to 0.41) • 24-months: 0.32 (0.23 to 0.41) • Statistically significant increase ($p \leq 0.01$) z-score (range) • Start of gender-affirming hormones: -0.27 ((-1.91 to 1.29) • 24-months: 0.02 (-2.1 to 1.35) • Statistically significant increase ($p \leq 0.05$) 	

Appendix F Quality appraisal checklists

Newcastle-Ottawa Quality Assessment Form for Cohort Studies

Note: A study can be given a maximum of one star for each numbered item within the Selection and Outcome categories. A maximum of two stars can be given for Comparability.

Selection

- 1) Representativeness of the exposed cohort
 - a) Truly representative (one star)
 - b) Somewhat representative (one star)
 - c) Selected group
 - d) No description of the derivation of the cohort
- 2) Selection of the non-exposed cohort
 - a) Drawn from the same community as the exposed cohort (one star)
 - b) Drawn from a different source
 - c) No description of the derivation of the non exposed cohort
- 3) Ascertainment of exposure
 - a) Secure record (e.g., surgical record) (one star)
 - b) Structured interview (one star)
 - c) Written self report
 - d) No description
 - e) Other
- 4) Demonstration that outcome of interest was not present at start of study
 - a) Yes (one star)
 - b) No

Comparability

- 1) Comparability of cohorts on the basis of the design or analysis controlled for confounders
 - a) The study controls for age, sex and marital status (one star)
 - b) Study controls for other factors (list) _____
(one star)
 - c) Cohorts are not comparable on the basis of the design or analysis controlled for confounders

Outcome

- 1) Assessment of outcome
 - a) Independent blind assessment (one star)
 - b) Record linkage (one star)
 - c) Self report
 - d) No description
 - e) Other
- 2) Was follow-up long enough for outcomes to occur
 - a) Yes (one star)
 - b) No

Indicate the median duration of follow-up and a brief rationale for the assessment above: _____
- 3) Adequacy of follow-up of cohorts
 - a) Complete follow up- all subject accounted for (one star)

- b) Subjects lost to follow up unlikely to introduce bias- number lost less than or equal to 20% or description of those lost suggested no different from those followed. (one star)
- c) Follow up rate less than 80% and no description of those lost
- d) No statement

Thresholds for converting the Newcastle-Ottawa scales to AHRQ standards (good, fair, and poor):

Good quality: 3 or 4 stars in selection domain AND 1 or 2 stars in comparability domain AND 2 or 3 stars in outcome/exposure domain

Fair quality: 2 stars in selection domain AND 1 or 2 stars in comparability domain AND 2 or 3 stars in outcome/exposure domain

Poor quality: 0 or 1 star in selection domain OR 0 stars in comparability domain OR 0 or 1 stars in outcome/exposure domain

Appendix G Grade profiles

Table 2: Question 1: For children and adolescents with gender dysphoria, what is the clinical effectiveness of treatment with gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? - Gender dysphoria

Study	QUALITY				Summary of findings				CERTAINTY	
	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect			IMPORTANCE
					Intervention	Comparator	Result	Result		
1 cohort study Lopez de Lara et al. 2020	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=23	None	T0 (baseline) = 57.1 (SD 4.1) T1 (12 months) = 14.7 (SD 3.2) Statistically significant improvement, p<0.001	Critical	VERY LOW	

Abbreviations: p: p-value; SD: standard deviation; UGDS: Utrecht Gender Dysphoria Scale

¹ Downgraded 1 level - the cohort study by Lopez de Lara et al. 2020 was assessed at high risk of bias (poor quality overall; lack of blinding and no control group)

Table 3: Question 1: For children and adolescents with gender dysphoria, what is the clinical effectiveness of treatment with gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? – Mental health

Study	QUALITY				Summary of findings				CERTAINTY	
	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect			IMPORTANCE
					Intervention	Comparator	Result	Result		
Impact on mental health (3 uncontrolled, prospective observational studies and 2 uncontrolled, retrospective observational studies)										
Change from baseline in mean depression score, measured using the BDI-II (duration of treatment 12 months). Higher scores indicate more severe depression.										

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events			
					Intervention	Comparator	Result	
1 cohort study Lopez de Lara et al. 2020	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=23	None	T0 (baseline) = 19.3 (SD 5.5) T1 (12 months) = 9.7 (SD 3.9) Statistically significant improvement, p<0.001	Critical VERY LOW
Change from baseline in mean depression score, measured using the CES-D-R (approximately 12-month follow-up). Higher scores indicate more severe depression.								
1 cohort study Achille et al. 2020	Serious limitations ²	Serious indirectness ³	No serious inconsistency	Not calculable	N=50	None	Wave 1 (baseline) = 21.4 Wave 3 (approx. 12 months) = 13.9 Statistically significant improvement (p<0.001)	Critical VERY LOW
Change from baseline in depression score, measured using the Patient Health Questionnaire Modified for Teens (PHQ 9_Modified for Teens) (approximately 12-month follow-up). Higher scores indicate more severe depression.								
1 cohort study Achille et al. 2020	Serious limitations ²	Serious indirectness ³	No serious inconsistency	Not calculable	N=50	None	Statistically significant reductions in mean score, p<0.001 Results presented diagrammatically, numerical results for mean score not reported	Critical VERY LOW
Change from baseline in depression symptoms, measured using the Quick Inventory of Depressive Symptoms (QIDS), self-reported (mean duration of gender-affirming hormone treatment 10.9 months). Higher scores indicate more severe depression.								
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=105	None	Baseline = 9.6 (SD 5.0) Follow-up = 7.4 (SD 4.5) No statistical analysis reported for the sub-group of participants receiving gender-affirming hormones	Critical VERY LOW
Change from baseline in depression symptoms, measured using the Quick Inventory of Depressive Symptoms (QIDS), clinician-reported (mean duration of gender-affirming hormone treatment 10.9 months). Higher scores indicate more severe depression.								
1 cohort study	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=106	None	Baseline = 5.9 (SD 4.1) Follow-up = 6.0 (SD 3.8)	Critical VERY LOW

QUALITY				Summary of findings				IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect		
					Intervention	Comparator	Result		
Kuper et al. 2020							No statistical analysis reported for the sub-group of participants who received gender-affirming hormones		
Need for treatment due to depression, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiata et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 54% (28/52) During real life phase 15% (8/52) Statistically significant reduction (p<0.001)	Critical	VERY LOW
Change from baseline in anxiety score, measured using the STAI-State subscale (duration of treatment 12 months). Higher scores indicate more severe anxiety.									
1 cohort study Lopez de Lara et al. 2020	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=23	None	T0 (baseline) = 33.3 (SD 9.1) T1 (12 months) = 16.8 (SD 8.1) Statistically significant improvement, p<0.001	Critical	VERY LOW
Change from baseline in anxiety score, measured using the STAI-Trait subscale (duration of treatment 12 months). Higher scores indicate more severe anxiety.									
1 cohort study Lopez de Lara et al. 2020	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=23	None	T0 (baseline) = 33.0 (SD 7.2) T1 (12 months) = 18.5 (SD 8.4) Statistically significant improvement, p<0.001	Critical	VERY LOW
Change from baseline in anxiety symptoms, measured using the SCARED questionnaire (mean duration of gender-affirming hormone treatment 10.9 months). Higher scores indicate more severe anxiety.									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=80	None	Baseline = 32.6 (SD 16.3) Follow-up = 28.4 (SD 15.9) No statistical analysis reported for the sub-group of participants	Critical	VERY LOW

QUALITY				Summary of findings				IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect		
					Intervention	Comparator	Result		
							who received gender-affirming hormones		
Change from baseline in panic symptoms, measured using specific questions from the SCARED questionnaire (mean duration of gender-affirming hormone treatment 10.9 months). Higher scores indicate more severe symptoms.									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=82	None	Baseline = 8.1 (SD 6.3) Follow-up = 7.1 (SD 6.5) No statistical analysis reported for the sub-group of participants who received gender-affirming hormones	Critical	VERY LOW
Change from baseline in generalised anxiety symptoms, measured using specific questions from the SCARED questionnaire (mean duration of gender-affirming hormone treatment was 10.9 months). Higher scores indicate more severe symptoms.									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=82	None	Baseline = 10.0 (SD 5.1) Follow-up = 8.8 (SD 5.0) No statistical analysis reported for the sub-group of participants who received gender-affirming hormones	Critical	VERY LOW
Change from baseline in social anxiety symptoms, measured using specific questions from the SCARED questionnaire (mean duration of gender-affirming hormone treatment was 10.9 months). Higher scores indicate more severe symptoms.									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=82	None	Baseline = 8.5 (SD 4.1) Follow-up = 7.7 (SD 4.2) No statistical analysis reported for the sub-group of participants who received gender-affirming hormones	Critical	VERY LOW
Change from baseline in separation anxiety symptoms, measured using specific questions from the SCARED questionnaire (mean duration of gender-affirming hormone treatment was 10.9 months). Higher scores indicate more severe symptoms.									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=81	None	Baseline = 3.5 (SD 3.0) Follow-up = 3.1 (SD 2.5) No statistical analysis reported for the sub-group of participants	Critical	VERY LOW

QUALITY				Summary of findings				IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect		
					Intervention	Comparator	Result		
Change from baseline in school avoidance, measured using specific questions from the SCARED questionnaire (mean duration of gender-affirming hormone treatment was 10.9 months). Higher scores indicate more severe symptoms.									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	No serious indirectness	No serious inconsistency	Not calculable	N=80	None	Baseline = 2.6 (SD 2.1) Follow-up = 2.0 (SD 2.0) No statistical analysis reported for the sub-group of participants who received gender-affirming hormones	Critical	VERY LOW
Need for treatment due to anxiety, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaltiala et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 48% (25/52) During real life phase 15% (8/52) Statistically significant reduction (p<0.001)	Critical	VERY LOW
Change from baseline in adjusted mean suicidality score, measured using the ASQ instrument (mean treatment duration 349 days). Higher scores indicate a greater degree of suicidality.									
1 cohort study Allen et al. 2019	Serious limitations ⁵	No serious indirectness	No serious inconsistency	Not calculable	N=39	None	T0 (baseline) = 1.11 (SE 0.22) T1 (final assessment) = 0.27 (SE 0.12) Statistically significant improvement in score from T0 to T1, p<0.001	Critical	VERY LOW
Change from baseline in percentage of participants with suicidal ideation, measured using the additional questions from the PHQ 9 Modified for Teens (approximately 12-month follow-up)									
1 cohort study Achille et al. 2020	Serious limitations ²	Serious indirectness ³	No serious inconsistency	Not calculable	N=50	None	Wave 1 (baseline) = 10% (5/50) Wave 3 (approx. 12 months) = 6% (3/50)	Critical	VERY LOW

QUALITY				Summary of findings				IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect Result		
					Intervention	Comparator			
Change from baseline in suicidal ideation (passive), information on which was collected by clinician, exact methods / tools not reported (mean duration of gender-affirming hormone treatment was 10.9 months)									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	Serious indirectness ⁶	No serious inconsistency	Not calculable	N=130	None	Lifetime = 81% (105 people) 1 month before initial assessment = 25% (33 people) Follow-up period = 38% (51 people) No statistical analysis reported	Critical	VERY LOW
Change from baseline in suicide attempts, information on which was collected by clinician, exact methods / tools not reported (mean duration of gender-affirming hormone treatment was 10.9 months)									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	Serious indirectness ⁶	No serious inconsistency	Not calculable	N=130	None	Lifetime = 15% (20 people) 3 months before initial assessment = 2% (3 people) Follow-up period = 5% (6 people) No statistical analysis reported	Critical	VERY LOW
Change from baseline in non-suicidal self-injury, information on which was collected by clinician, exact methods / tools not reported (mean duration of gender-affirming hormone treatment was 10.9 months)									
1 cohort study Kuper et al. 2020	Serious limitations ⁴	Serious indirectness ⁶	No serious inconsistency	Not calculable	N=130	None	Lifetime = 52% (68 people) 3 months before initial assessment = 10% (13 people) Follow-up period = 17% (23 people) No statistical analysis reported	Critical	VERY LOW
Need for treatment due to suicidality / self-harm, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiata et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 35% (18/52) During real life phase	Critical	VERY LOW

QUALITY					Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect		
					Intervention	Comparator	Result		
							4% (2/52) Statistically significant reduction (p<0.001)		
Need for mental health treatment, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiala et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 50% (26/52) During real life phase 46% (24/51) No statistically significant difference (p= 0.77)	Critical	VERY LOW
Need for treatment due to conduct problems / antisocial, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiala et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 14% (7/52) During real life phase 6% (3/52) No statistically significant difference (p= 0.18)	Critical	VERY LOW
Need for treatment due to psychotic symptoms or psychosis, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiala et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 2% (1/52) During real life phase 4% (2/52) No statistically significant difference (p= 0.56)	Critical	VERY LOW
Need for treatment due to substance abuse, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									

QUALITY					Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect		
					Intervention	Comparator	Result		
1 cohort study Kaitiata et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 4% (2/52) During real life phase 2% (1/52) No statistically significant difference (p= 0.56)	Critical	VERY LOW
Need for treatment due to autism, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiata et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 12% (6/52) During real life phase 6% (3/52) No statistically significant difference (p= 0.30)	Critical	VERY LOW
Need for treatment due to ADHD, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiata et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 10% (5/52) During real life phase 2% (1/52) No statistically significant difference (p= 0.09)	Critical	VERY LOW
Need for treatment due to eating disorder, during and before gender identity assessment, and during real life phase (approximately 12 months follow-up)									
1 cohort study Kaitiata et al. 2020	Serious limitations ⁷	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During and before gender identity assessment 2% (1/52)	Critical	VERY LOW

QUALITY					Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of events		Effect		
					Intervention	Comparator	Result		
							During real life phase 2% (1/52) No statistically significant difference (p=1.0)		

Abbreviations: ADHD: attention deficit hyperactivity disorder; ASQ: Ask Suicide-Screening Questions; CESD-R: Center for Epidemiologic Studies Depression Scale; BDI-II: Beck Depression Inventory II (BDI-II); p: p-value; PHQ 9_Modified for Teens: Patient Health Questionnaire Modified for Teens; SCARED: Screen for Child Anxiety Related Emotional Disorders; SD: standard deviation; STAI: State-Trait Anxiety Inventory

- 1 Downgraded 1 level - the cohort study by Lopez de Lara et al. (2020) was assessed at high risk of bias (poor quality; lack of blinding and no control group).
- 2 Downgraded 1 level - the cohort study by Achille et al (2020) was assessed at high risk of bias (poor quality; lack of blinding, no control group and high number of participants lost to follow-up).
- 3 Serious indirectness in Achille 2020- Outcome reported for full study cohort, of whom 30% were taking no treatment or puberty suppression alone at follow-up. Results for people taking gender-affirming hormones not reported separately.⁴ Downgraded 1 level - the cohort study by Kuper et al. (2020) was assessed at high risk of bias (poor quality).
- 5 Downgraded 1 level - the cohort study by Allen et al. (2019) was assessed at high risk of bias (poor quality; lack of blinding and no control group).
- 6 Serious indirectness in Kuper et al. 2020- Outcome reported for full study cohort, of whom approximately 17% received puberty suppression alone and did not receive gender-affirming hormones
- 7 Downgraded 1 level - the cohort study by Kalliala et al. (2020) was assessed at high risk of bias (poor quality; lack of blinding and no control group).

Table 4: Question 1: For children and adolescents with gender dysphoria, what is the clinical effectiveness of treatment with gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? – Quality of life

QUALITY							Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect				
					Intervention	Comparator	Result				
Impact on quality of life (1 uncontrolled, prospective observational study and 1 uncontrolled, retrospective observational study)											

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY	
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients	Effect			
					Intervention	Comparator	Result		
Change from baseline in mean quality of life score, measured using the QLES-Q-SF) (approximately 12-month follow-up). Higher scores indicated better quality of life.									
1 cohort study Achille et al. 2020	Serious limitations ¹	Serious indirectness ²	No serious inconsistency	Not calculable	N=50	None	Numerical improvements in mean score reported from wave 1 (baseline) to wave 3 (approx. 12 months), but difference not statistically significant (p = 0.085) Results presented diagrammatically, numerical results for mean score not reported	Critical	VERY LOW
Change from baseline in adjusted mean well-being score, measured using the GWBS of the Pediatric Quality of Life Inventory (mean treatment duration 349 days). Higher scores indicated better well-being.									
1 cohort study Allen et al. 2019	Serious limitations ³	No serious indirectness	No serious inconsistency	Not calculable	N=39	None	T0 (baseline) = 61.70 (SE 2.43) T1 (final assessment) = 70.23 (SE 2.15) Statistically significant improvement in well-being score, p<0.002	Critical	VERY LOW

Abbreviations: GWBS: General Well-Being Scale; p: p-value; QLES-Q-SF: Quality of Life Enjoyment and Satisfaction Questionnaire; SE: standard error

¹ Downgraded 1 level - the cohort study by Achille et al (2020) was assessed at high risk of bias (poor quality; lack of blinding, no control group and high number of participants lost to follow-up).
² Serious indirectness in Achille et al. 2020 - Outcome reported for full study cohort, of whom 30% were taking no treatment or puberty suppression alone at follow-up. Results for people taking gender-affirming hormones not reported separately.
³ Downgraded 1 level - the cohort study by Allen et al. (2019) was assessed at high risk of bias (poor quality; lack of blinding and no control group).

Table 5: Question 1: For children and adolescents with gender dysphoria, what is the clinical effectiveness of treatment with gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? – Body image

QUALITY			Summary of findings			IMPORTANCE	CERTAINTY
---------	--	--	---------------------	--	--	------------	-----------

Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect
					Intervention	Comparator	
Impact on body image (1 uncontrolled, prospective observational study)							
Change from baseline in mean body image, measured using the BIS (mean duration of gender-affirming hormone treatment was 10.9 months). Higher scores represent a higher degree of body dissatisfaction.							
1 cohort study Kuper et al. 2020	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=86	None	Baseline = 70.7 (SD 15.2) Follow-up = 51.4 (SD 18.3) No statistical analysis reported for the sub-group of participants who received gender-affirming hormones Important VERY LOW

Abbreviations: BIS: Body Image Scale; p: p-value; SD: standard deviation

¹ Downgraded 1 level - the cohort study by Kuper et al. (2020) was assessed at high risk of bias (poor quality; lack of blinding, no control group and high number of participants lost to follow-up).

Table 6: Question 1: For children and adolescents with gender dysphoria, what is the clinical effectiveness of treatment with gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? – Psychological impact

Study	QUALITY				Summary of findings			CERTAINTY
	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect	
					Intervention	Comparator		
Psychosocial Impact (1 uncontrolled, prospective observational study and 1 uncontrolled, retrospective observational study)								
Change from baseline in family functioning, measured using the Family APGAR test. Higher scores suggest more family dysfunction.								
1 cohort study Lopez de Lara et al. 2020	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=23	None	T0 (baseline) = 17.9 T1 (12 months) = 18.0 No statistical analysis reported	Important VERY LOW
Change from baseline in mean patient strengths and difficulties score, measured using the SDQ, Spanish Version (total difficulties score) (duration of treatment 12 months). Higher scores suggest the presence of a behavioural disorder.								
1 cohort study	Serious limitations ¹	No serious indirectness	No serious inconsistency	Not calculable	N=23	None	T0 (baseline) = 14.7 (SD 3.3) T1 (12 months) = 10.3 (SD 2.9)	Important VERY LOW

This document was prepared in October 2020

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect		
Study				Intervention	Comparator	Result (95% CI)		
Lopez de Lara et al. 2020						Statistically significant improvement p<0.001		
Functioning in adolescent development: Living with parent(s)/ guardians² (outcome reported for the approximately 12-month period after starting gender-affirming hormones; referred to as the 'real-life phase' in Finland). Not living with parent(s) or guardian in your early 20s is a marker of age-appropriate functioning in Finnish culture.								
1 cohort study Kaitiainen et al. 2020	Serious limitations ³	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During gender identity assessment = 73% (38/52) During real life phase = 40% (21/50) Statistically significant reduction (p=0.001)	Important VERY LOW
Functioning in adolescent development: Normative peer contacts⁴ (outcome reported for the approximately 12-month period after starting gender-affirming hormones; referred to as the 'real-life phase' in Finland)								
1 cohort study Kaitiainen et al. 2020	Serious limitations ³	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During gender identity assessment = 89% (46/52) During real life phase = 81% (42/52) Statistically significant reduction (p<0.001)	Important VERY LOW
Functioning in adolescent development: Progresses normatively in school/ work⁵ (outcome reported for the approximately 12-month period after starting gender-affirming hormones; referred to as the 'real-life phase' in Finland)								
1 cohort study Kaitiainen et al. 2020	Serious limitations ³	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During gender identity assessment = 64% (33/52) During real life phase = 60% (31/52) No statistically significant difference (p=0.69)	Important VERY LOW
Functioning in adolescent development: Has been dating or had steady relationships⁶ (outcome reported for the approximately 12-month period after starting gender-affirming hormones; referred to as the 'real-life phase' in Finland)								
1 cohort study	Serious limitations ³	No serious indirectness	No serious inconsistency	Not calculable	N=52	None	During gender identity assessment = 62% (32/50)	Important VERY LOW

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	No of patients		Effect		
				Intervention	Comparator		Result (95% CI)	
Kaltiala et al. 2020						During real life phase = 58% (30/52) No statistically significant difference (p=0.51)		
Functioning in adolescent development: Is age-appropriately able to deal with matters outside of the home⁷ (outcome reported for the approximately 12-month period after starting gender-affirming hormones; referred to as the 'real-life phase' in Finland)								
1 cohort study Kaltiala et al. 2020	Serious limitations ²	No serious indirectness	No serious inconsistency	N=52	None	During gender identity assessment = 81% (42/52) During real life phase = 81% (42/52) No statistically significant difference (p=1.00)	Important	VERY LOW

Abbreviations: APGAR: Adaptability, Partnership, Growth, Affection and Resolve; p: p-value; SD: standard deviation; SDQ: Strengths and Difficulties Questionnaire

1 Downgraded 1 level - the cohort study by Lopez de Lara et al. (2020) was assessed at high risk of bias (poor quality; lack of blinding and no control group).

2 Living arrangements were classified as (1) living with at least one parent/guardian, (2) living in a boarding school, with an adult relative, in some form of supported accommodation or the like, where supervision and guidance by a responsible adult is provided, (3) independently alone or in a shared household with a peer, (4) with a romantic partner. In the analyses dichotomised living arrangements as (a) parent(s)/guardian(s) vs. in other arrangements.

3 Downgraded 1 level - the cohort study by Kaltiala et al. (2020) was assessed at high risk of bias (poor quality; lack of blinding and no control group).

4 Peer relationships were classified as: (1) socialises with friends in leisure time, outside of activities supervised by adults, (2) socialises with peers only at school or in the context of rehabilitative activity, (3) spends time close to peers, for example in school or rehabilitative activity, but does not connect with them, (4) does not meet peers at all. In the analyses, peer relationships during (a) gender identity assessment and (b) the real-life phase were dichotomized to age-appropriate (normative) (1) vs. restricted or lacking (2–4).

5 School/work participation was classified as (1) age appropriate participation in mainstream curriculum, progresses without difficulties, (2) participates in mainstream curriculum with difficulty, (3) participates in rehabilitative educational or work activity, (4) not involved in education and working life. Age-appropriate participation during (1) was recorded if the adolescent attended mainstream secondary education or upper secondary education at a regular rate (a class per year in comprehensive school; has not changed more than once between tracks in upper secondary education) or had proceeded to work life after completing vocational education. Participation with difficulty (2) was recorded if the adolescent was enrolled in mainstream education but had to repeat a class, studied with special arrangements (for example, in a special small group), or followed some form of adjusted curriculum. In the analyses, school/work life during (a) gender identity assessment and (b) real-life phase was dichotomised to normative (1) vs. any other (2, 3 or 4).

6 Romantic involvement was recorded (1) has or has had a dating or steady relationship, not only online, (2) has had a romantic relationship only online, (3) has not had dating or steady relationships. In the analyses we compared has or has had (1) vs. has not had (2,3) a dating or steady relationship during (a) gender identity assessment and (b) real-life phase. Sexual history was recorded in more detail in case histories during gender identity assessment, and for this period we also collected the experiences of (French) kissing (yes/no), intercourse (yes/no) and experience of any genitally intimate contact with a partner (petting under clothes or naked, intercourse, oral sex) (yes/no).

7 In recording age-appropriate competence in managing everyday matters it was expected that early adolescents (up to 14 years) would be able, for example, to do shopping and travel alone on local public transport, and to help with household duties assigned by their parents. Middle adolescents (15–17 years) were further assumed, for example, to be able make telephone calls in matters important to them (for example, when seeking a summer job), to deal with school-related issues with school personnel without parental participation, to select and start new hobbies independently and to fulfil their role in summer jobs and in similar responsibilities of young people. Late adolescents (18 years and over), legally adults, were expected to have, in addition to the above, competence to talk to authorities such as professionals in health and social services, employment or educational institutions, to deal with banks or health insurance, to manage their financial issues and to manage their housekeeping if they chose to move to live independently of parents/guardians. Competence in managing everyday matters was recorded as follows: (1) the adolescent is able to cope age appropriately outside home, (2) the adolescent needs support in age-appropriate matters outside home but functions age-appropriately in the home (manages her/his own hygiene, clothing and nutrition, participates in (younger subjects) or takes responsibility for (older subjects) housekeeping) and (3) the adolescent's functioning is inadequate both at home and outside home. For the analyses, participants were determined to be able to age-appropriately able cope with matters outside of the home (1) vs. not (2,3).

Table 7: Question 2: For children and adolescents with gender dysphoria, what is the short-term and long-term safety of gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? – Bone density

Study	QUALITY					Summary of findings			IMPORTANCE	CERTAINTY
	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect			
					Intervention	Comparator		Result (95% CI)		
Lumbar spine bone mineral apparent density (BMAD) (2 uncontrolled, retrospective observational studies)										
Change from start of gender-affirming hormones to age 22 years in lumbar spine BMAD in transfemales										
1 cohort study Klink et al. 2015	Serious limitations ¹	Serious indirectness ²	Not applicable	Not calculable	N=13 (Mean) N=14 (z-score)	None	Mean (SD), g/m ³ Start of gender-affirming hormones: 0.22 (0.02) Age 22 years: 0.23 (0.03) P=0.003 z-score (SD) Start of gender-affirming hormones: -0.90 (0.80) Age 22 years: -0.78 (1.03) No statistically significant difference	Important	VERY LOW	
Change from baseline in lumbar spine BMAD in transfemales with a bone age less than 15 years (young); 24 months follow-up)										
1 cohort study Vlot et al. 2017	Serious limitations ³	No serious indirectness	Not applicable	Not calculable	N=15	None	Median (range), g/m ³ Start of gender-affirming hormones (CO): 0.20 (0.18 to 0.24)	Important	VERY LOW	

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY	
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients				Effect
					Intervention	Comparator			
1 cohort study Vlot et al. 2017	Serious limitations ³	No serious indirectness	Not applicable	Not calculable	N=10	None	Median (range), g/m ³ C0: 0.30 (0.22 to 0.35) C24: 0.33 (0.23 to 0.37) Statistically significant increase (p≤0.01) z-score (range) C0: -0.37 (-2.28 to 0.47) C24: -0.37 (-2.03 to 0.85) Statistically significant increase (p≤0.01)	Important	VERY LOW
Change from baseline in femoral neck BMAD in transmales with a bone age of 14 years or more ('old'; 24 months follow-up)									
1 cohort study Vlot et al. 2017	Serious limitations ³	No serious indirectness	Not applicable	Not calculable	N=23	None	Median (range), g/m ³ C0: 0.30 (0.23 to 0.41) C24: 0.32 (0.23 to 0.41) Statistically significant increase (p≤0.01) z-score (range) C0: -0.27 (-1.91 to 1.29) C24: 0.02 (-2.1 to 1.35) Statistically significant increase (p≤0.05)	Important	VERY LOW
Change in lumbar spine BMD (2 uncontrolled, retrospective observational studies)									
Change from start of gender-affirming hormones to age 22 years in lumbar spine BMD in transfemales									
1 cohort study Klink et al. 2015	Serious limitations ¹	Serious indirectness ²	Not applicable	Not calculable	N=15 (Mean) N=13 (z-score)	None	Mean (SD), g/m ² Start of gender-affirming hormones: 0.84 (0.11) Age 22 years: 0.93 (0.10) P<0.001 z-score (SD)	Important	VERY LOW

QUALITY					Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect		
					Intervention	Comparator			
Change from start of gender-affirming hormones to age 22 years in lumbar spine BMD in transmales									
1 cohort study Klink et al. 2015	Serious limitations ¹	Serious indirectness ²	Not applicable	Not calculable	N=19 (Mean and z-score)	None	Mean (SD), g/m ² Start of gender-affirming hormones: 0.91 (0.10) Age 22 years: 0.99 (0.13) P<0.001 z-score (SD) Start of gender-affirming hormones: -0.72 (0.99) Age 22 years: -0.33 (1.12) No statistically significant difference	Important	VERY LOW
Change from start of testosterone treatment in lumbar spine BMD in transmen (follow-up 6 to 24 months)									
1 cohort study Stoffers et al. 2019	Serious limitations ⁴	No serious indirectness	Not applicable	Not calculable	N=62 (T0 and T6) N=37 (T12) N=15 (T24)	None	Mean (SD), g/cm ² T0: 0.90 (0.11) T6: 0.94 (0.10) T12: 0.95 (0.09) T24: 0.95 (0.11) No statistically significant difference from T0 to any timepoint z-score (SD) T0: -0.81 (1.02) T6: -0.67 (0.95) T12: -0.66 (0.81) T24: -0.74 (1.17)	Important	VERY LOW

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients	Effect		
					Intervention	Comparator	Result (95% CI)	
Stoffers et al. 2019					N=37 (T12) N=15 (T24)		T6: 0.84 (0.11) T12: 0.82 (0.08) T24: 0.85 (0.11) No statistically significant difference from T0 to any timepoint z-score (SD) T0: -0.97 (0.79) T6: -0.54 (0.96) T12: -0.80 (0.69) T24: -0.31 (0.84) No statistically significant difference from T0 to any timepoint	
Change from start of testosterone treatment in left femoral neck (hip) BMD in transmales (follow-up 6 to 24 months)								
1 cohort study Stoffers et al. 2019	Serious limitations ⁴	No serious indirectness	Not applicable	Not calculable	N=62 (T0 and T6) N=37 (T12) N=15 (T24)	None	Mean (SD), g/cm ² T0: 0.76 (0.09) T6: 0.83 (0.12) T12: 0.81 (0.08) T24: 0.86 (0.09) No statistically significant difference from T0 to any timepoint z-score (SD) T0: -1.07 (0.85) T6: -0.62 (1.12) T12: -0.93 (0.63) T24: -0.20 (0.70) No statistically significant difference from T0 to any timepoint	Important VERY LOW

Abbreviations: BMAD: bone mineral apparent density; BMD: bone mineral density; g: grams; m: metre; SD: standard deviation

- 1 Downgraded 1 level - the cohort study by Klink et al. (2015) was assessed as at high risk of bias (poor quality overall; lack of blinding, no control group and high number of participants lost to follow-up)
- 2 Outcomes reported after gender reassignment surgery and not after gender-affirming hormones alone. Unclear whether observed changes are due to hormones or surgery
- 3 Downgraded 1 level - the cohort study by Vlot et al. (2017) was assessed as at high risk of bias (poor quality overall; lack of blinding and no control)
- 4 Downgraded 1 level - the cohort study by Stoffers et al. (2019) was assessed as at high risk of bias (poor quality overall; lack of blinding and no control group)

Table 8: Question 2: For children and adolescents with gender dysphoria, what is the short-term and long-term safety of gender-affirming hormones compared with one or a combination of psychological support, social transitioning to the desired gender or no intervention? – Cardiovascular risk factors

Study	QUALITY				Summary of findings			CERTAINTY
	Risk of bias	Indirectness	Inconsistency	Imprecision	No. of patients	Effect	IMPORTANCE	
Change in body mass index (1 uncontrolled, retrospective observational study)								
Change from start of gender-affirming hormones to age 22 years in BMI in transfemales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=71	None	Mean change (95% CI) +1.9 (0.6 to 3.2) Statistically significant increase (p<0.005) Mean BMI at 22 years (95% CI): 23.2 (21.6 to 24.8)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in BMI in transmales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI) +1.4 (0.8 to 2.0) Statistically significant increase (p<0.005) Mean BMI at 22 years (95% CI): 23.9 (23.0 to 24.7)	Important VERY LOW

QUALITY					Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients		Effect		
					Intervention	Comparator	Result (95% CI)		
Obesity rates at age 22 years in transfemales who started gender-affirming hormones as adolescents (1 uncontrolled, retrospective observational study)									
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=71	None	At 22 years, 9.9% of transfemales were obese, compared with 3.0% in reference cisgender population No statistically analysis reported	Important	VERY LOW
Obesity rates at age 22 years in transfemales who started gender-affirming hormones as adolescents (1 uncontrolled, retrospective observational study)									
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	At 22 years, 6.6% of transfemales were obese, compared with 2.2% in reference cisgender population No statistically analysis reported	Important	VERY LOW
Change in blood pressure (1 uncontrolled, retrospective observational study)									
Change in blood pressure (1 uncontrolled, retrospective observational study)									
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=71	None	Mean change (95% CI) -3 (-8 to 2) No statistically significant difference Mean SBP at 22 years (95% CI): 117 (113 to 122)	Important	VERY LOW
Change from start of gender-affirming hormones to age 22 years in diastolic blood pressure (DBP) in transfemales									

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients			
					Intervention	Comparator		
							No statistically significant difference Mean glucose level at 22 years (95% CI): 5.0 (4.8 to 5.1)	
Change from start of gender-affirming hormones to age 22 years in insulin level (mU/L) in transfemales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=71	None	Mean change (95% CI) +2.7 (-1.7 to 7.1) No statistically significant difference Mean insulin level at 22 years (95% CI): 13.0 (8.4 to 17.6)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in insulin resistance (HOMA-IR) in transfemales. Higher scores indicate more insulin resistance.								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=71	None	Mean change (95% CI) +0.7 (-0.2 to 1.5) No statistically significant difference Mean HOMA-IR at 22 years (95% CI): 2.9 (1.9 to 3.9)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in glucose level (mmol/L) in transmales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI) 0.0 (-0.2 to 0.2) No statistically significant difference	Important VERY LOW

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients			
					Intervention	Comparator		
Change from start of gender-affirming hormones to age 22 years in insulin level (mU/L) in transmales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI) -2.1 (-3.9 to -0.3) Statistically significant decrease (p<0.05) Mean insulin level at 22 years (95% CI): 8.6 (6.9 to 10.2)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in insulin resistance (HOMA-IR) in transmales. Higher scores indicate more insulin resistance.								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI): -0.5 (-1.0 to -0.1) Statistically significant decrease (p<0.05) Mean HOMA-IR at 22 years (95% CI): 1.8 (1.4 to 2.2)	Important VERY LOW
Change from start of testosterone in HbA1c in transmales (up to 24 months follow-up)								
1 cohort study Stoffers et al. 2019	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N= Not reported	None	No statistically significant change from start of testosterone treatment Numerical results, follow-up duration and further details of statistical analysis not reported.	Important VERY LOW

QUALITY				Summary of findings			IMPORTANCE	CERTAINTY
Study	Risk of bias	Indirectness	Inconsistency	Imprecision	No of patients			
					Intervention	Comparator	Result (95% CI)	
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=71	None	Mean change (95% CI): +0.2 (0.0 to 0.5) Statistically significant increase (p<0.05) Mean triglycerides at 22 years (95% CI): 1.1 (0.9 to 1.4)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in total cholesterol (mmol/L) in transmales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI): +0.4 (0.2 to 0.6) Statistically significant increase (p<0.001) Mean total cholesterol at 22 years (95% CI): 4.6 (4.3 to 4.8)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in HDL cholesterol (mmol/L) in transmales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI): -0.3 (-0.4 to -0.2) Statistically significant decrease (p<0.001) Mean HDL cholesterol at 22 years (95% CI): 1.3 (1.2 to 1.3)	Important VERY LOW
Change from start of gender-affirming hormones to age 22 years in LDL cholesterol (mmol/L) in transmales								
1 cohort study Klaver et al. 2020	Serious limitations ¹	No serious indirectness	Not applicable	Not calculable	N=121	None	Mean change (95% CI): +0.4 (0.2 to 0.6) Statistically significant increase (p<0.001)	Important VERY LOW